

Putting the Precision in Precision Cosmology: How accurate should your data covariance matrix be?

Andy Taylor^{1*}, Benjamin Joachimi¹, Thomas Kitching^{1,2}

1. *Scottish Universities Physics Alliance (SUPA), Institute for Astronomy, School of Physics and Astronomy, University of Edinburgh, Royal Observatory, Blackford Hill, Edinburgh, EH9 3HJ, U.K.*

2. *Mullard Space Science Laboratory, University College London, Holmbury St. Mary, Surrey, RH5 6NT, U.K.*

19 December 2012

ABSTRACT

Cosmological parameter estimation requires that the likelihood function of the data is accurately known. Assuming that cosmological large-scale structure power spectra data are multivariate Gaussian-distributed, we show the accuracy of parameter estimation is limited by the accuracy of the inverse data covariance matrix – the *precision matrix*. If the data covariance and precision matrices are estimated by sampling independent realisations of the data, their statistical properties are described by the Wishart and Inverse-Wishart distributions, respectively. Independent of any details of the survey, we show that the fractional error on a parameter variance, or a Figure-of-Merit, is equal to the fractional variance of the precision matrix. In addition, for the only unbiased estimator of the precision matrix, we find that the fractional accuracy of the parameter error depends only on the difference between the number of independent realisations and the number of data points, and so can easily diverge. For a 5% error on a parameter error and $N_D \ll 10^2$ data-points, a minimum of 200 realisations of the survey are needed, with 10% accuracy in the data covariance. If the number of data-points $N_D \gg 10^2$ we need $N_S > N_D$ realisations and a fractional accuracy of $< \sqrt{2/N_D}$ in the data covariance. As the number of power spectra data points grows to $N_D > 10^4 - 10^6$ this approach will be problematic. We discuss possible ways to relax these conditions: improved theoretical modelling; shrinkage methods; data-compression; simulation and data resampling methods.

Key words: Cosmology: theory - large-scale structure of Universe, methods: statistical analysis

1 INTRODUCTION

A central part of modern cosmology is the measurement of the parameters that characterise cosmological models of the Universe. These can be the set that constitutes the Standard Cosmological Model ($\Omega_m, \Omega_b, \Omega_\Lambda, H_0, \sigma_8, n_s, \tau$), or an extended set that characterise, for example, more complex dark energy models (see e.g., Copeland, Sami & Tsujikawa, 2006; Amendola et al., 2012, for reviews), deviations from Einstein gravity (e.g., Clifton, Ferreira, Padillo, Skordis, 2012; Amendola et al., 2012 for recent reviews), more detail about the inflationary epoch (e.g., Amendola et al, 2012), isocurvature density and velocity modes (e.g., Bucher, Moodley & Turok, 2001), or massive neutrinos and their abundance (e.g., Bird, Viel & Haehnelt, 2012, and references therein). Furthermore, if we want to differentiate between theoretical models in a Bayesian framework,

as well as estimate their parameter value, we also need to accurately integrate over the model parameter-space (e.g., Trotta, 2007; Liddle, Mukherjee, Parkinson, 2006; Taylor & Kitching, 2010).

To carry these tasks out we need both accurate theoretical predictions of the physical properties of the model to compare to the data, and sufficiently accurate models of their statistical properties. Ideally, we would like to be able to accurately predict the full multivariate probability distribution of the data for each model. If, as is commonly assumed, the data can be modelled as a multivariate Gaussian distribution, all of the statistical properties of the model reside in the mean and covariance of the model. Attention has been focussed on the accuracy of the predictions of the mean value – e.g., the model power spectra – and the effect of biases or errors in the mean (e.g., Huterer & Takada, 2005; Huterer, et al., 2006; Taylor et al., 2007). But to fully specify the distribution of the data we also need accurate

* ant@roe.ac.uk

predictions of the data covariance matrix and the inverse of the data covariance – the *precision matrix*.

If we assume that the mean is well-known, the accuracy of the probability distribution of the data, and hence the likelihood function in parameter-space, is determined by the accuracy of the precision matrix. However, as yet there is no unique approach to estimating the data covariance matrix since this may depend on the details of what is known about the data, and even less attention paid to the estimation of the precision matrix.

The data covariance matrix can be estimated in a number of ways: direct calculation of a theoretical model; estimate the sample covariance from an ensemble of simulations of the data; or estimate the sample covariance matrix from the data itself. If we know the data covariance from the theoretical model accurately enough, there is no statistical uncertainty, and the precision matrix can be accurately estimated. But if the data covariance matrix must be sampled from an ensemble of simulations, or the data itself, there will be statistical uncertainty in the sample covariance. If we assume the underlying data is Gaussian-distributed and the samples are independent and drawn from the same distribution, the probability distribution of the sample covariance matrix is known, and was first derived by Wishart (1928; see also e.g., Press, 1982). To fully specify the model distribution of the data we also require the precision matrix. The distribution of the precision matrix, the Inverse-Wishart distribution (e.g., Press, 1982), has significantly different properties from distribution of the covariance matrix. The Wishart distributions has previously been discussed in cosmology as the distribution of Cosmic Microwave Background (CMB) temperature and polarisation power spectra (Percival & Brown, 2006), while the Inverse-Wishart distribution has been used as a prior for Bayesian estimates of the CMB temperature and polarisation power spectra (Eriksen & Wehus, 2009), for Gibbs sampling (Larson et al., 2007), and to test Pseudo-Cl methods (Hamimeche & Lewis, 2009).

In this paper we develop a new framework to estimate the statistical error on the data covariance and precision matrix (Section 3). We illustrate these effects on simulated data (Section 4) and discuss the implications for imminent and future large-scale structure surveys in cosmology. These effects are propagated into the accuracy of parameter errors, and the parameter covariance around the peak of the likelihood surface (Section 5). Since many experiments use the 2-parameter Figure-of-Merit (FoM) as a target measure for survey design, we also discuss the accuracy of an arbitrary FoM (Section 6). Given a prescribed accuracy for the parameter covariance matrix, or a FoM, we show how accurate the precision matrix and data covariance matrix must be. Finally, we discuss ways in which we avoid these bounds by improved theoretical modelling of the data covariance, rapid simulation production, or using data compression and shrinkage methods (Section 7). We begin by reviewing parameter estimation and the role of the precision matrix.

2 PARAMETER ESTIMATION

To begin with we shall assume that the cosmological parameters, θ , being measured are estimated from maximising a posterior parameter distribution, $p(\theta|\mathbf{D}, \mathcal{M})$, given a

dataset, \mathbf{D} , and some theoretical model, \mathcal{M} (see, e.g., Sivia, 1996). From Bayes Theorem,

$$p(\theta|\mathbf{D}, \mathcal{M}) = \frac{L(\mathbf{D}|\theta, \mathcal{M})\pi(\theta|\mathcal{M})}{E(\mathbf{D}|\mathcal{M})}, \quad (1)$$

we can determine the posterior parameter distribution from the likelihood function for the data, $L(\mathbf{D}|\theta, \mathcal{M})$, predicted by the model, a prior, $\pi(\theta|\mathcal{M})$, which is the probability distribution of the parameters before the data is analysed, and normalised by the evidence, $E(\mathbf{D}|\mathcal{M})$, which marginalises over the likelihood and prior in parameter-space. If we restrict our study to parameter estimation for a given model, we can ignore this term. We shall assume the prior on the parameters is flat.

If we model the data distribution as a multivariate Gaussian, then the likelihood function can be written

$$L(\mathbf{D}|\mu, \mathbf{M}, \mathcal{M}) = \frac{1}{(2\pi)^{N_D/2} \sqrt{|\mathbf{M}|}} \exp -\frac{1}{2} \text{Tr } \mathbf{W} \Psi, \quad (2)$$

where

$$\mathbf{W} = \Delta \mathbf{D} \Delta^t, \quad (3)$$

a superscript, t , indicates a transpose,

$$\Delta \mathbf{D} = \mathbf{D} - \langle \mathbf{D} \rangle \quad (4)$$

is the variation in the data-vector, $\mu = \langle \mathbf{D} \rangle$ is the mean of the data and N_D is the length of the data-vector. The data covariance matrix is given by

$$\mathbf{M} = \langle \mathbf{W} \rangle = \langle \Delta \mathbf{D} \Delta^t \rangle. \quad (5)$$

We define $|\mathbf{M}| = \det \mathbf{M}$ as the determinant. Comparing with a multivariate Gaussian we see that the matrix, Ψ , is the inverse of the data covariance matrix;

$$\Psi = \mathbf{M}^{-1}. \quad (6)$$

As we shall find this matrix is central to our analysis, we shall define the inverse data covariance as the *precision matrix*. The model dependence on cosmological model parameters, θ , may lie in either the mean, $\mu = \mu(\theta)$, or the data covariance matrix, $\mathbf{M} = \mathbf{M}(\theta)$, or both. Throughout we shall assume that the cosmological parameter dependence lies only in the mean. In Appendix A we describe the data vectors commonly used in cosmological large-scale structure analysis: galaxy redshift surveys, cosmic microwave background experiments and weak lensing surveys. Throughout, we shall assume that the data is a set of power spectra estimated from the data, although of course our results hold for correlation functions and are general to Gaussian-distributed data.

3 COVARIANCE AND PRECISION

3.1 Data covariance matrix

If we have a physical model for the covariance matrix, we would choose to use this. However, the statistical properties of the data may be poorly understood, for example the nonlinear regime for galaxy redshift surveys and weak lensing, and galaxy bias in redshift surveys, or the data may have been processed in ways which are not straightforward to model analytically, e.g., in CMB data where long-wave

variations in the time-ordered-data may have to be removed via polynomial fits, which can alter the statistical properties. In these cases we use an ensemble of simulations to estimate the sample data covariance matrix. In surveys where we do not know how to accurately simulate the data, we can use the data itself to estimate the data covariance. We return to this issue in Section 7.

If we generate N_S independent realisations of the data, \mathbf{D}_α , where each realisation is labelled by a Greek index, α, β, \dots , and adopt a convention of labelling the data-vector so that the Roman indices, i, j, \dots , indicates the wavenumber, ℓ , or wavevector, \mathbf{k} , and redshifts z, z', \dots , the data-vector averaged over the realisations is

$$\overline{\mathbf{D}} = \frac{1}{N_S} \sum_{\alpha} \mathbf{D}_{\alpha}. \quad (7)$$

The expectation value of the data-vector is

$$\langle \mathbf{D}_{\alpha} \rangle = \langle \overline{\mathbf{D}} \rangle = \boldsymbol{\mu}. \quad (8)$$

For independent and identically distributed realisations, and where we can use a symmetry or binning of the data to average over N_{modes} with the same mean value, the accuracy of the estimate of the mean data-vector will scale as

$$\sigma(\overline{\mathbf{D}}) = \sqrt{\frac{1}{N_S N_{\text{modes}}}} \boldsymbol{\mu}. \quad (9)$$

An unbiased estimator for the data covariance matrix is the sample data covariance, $\widehat{\mathbf{M}}$, from an ensemble of N_S independent and identically distributed realisations;

$$\widehat{\mathbf{M}} = \frac{1}{\nu} \sum_{\alpha} \Delta \mathbf{D}_{\alpha} \Delta \mathbf{D}_{\alpha}^t, \quad (10)$$

where

$$\Delta \mathbf{D}_{\alpha} = \mathbf{D}_{\alpha} - \overline{\mathbf{D}}_{\alpha} \quad (11)$$

is the variation in the data for each realisation, and ν is the number of degrees-of-freedom in the ensemble. If the estimated mean of the data-vector is known to be the expected mean, $\overline{\mathbf{D}}_{\alpha} = \langle \mathbf{D} \rangle$ then

$$\nu = N_S. \quad (12)$$

However, if the mean is estimated from the data itself, we reduce the number degrees-of-freedom by one, so that

$$\nu = N_S - 1. \quad (13)$$

3.2 The Wishart distribution

The statistical properties of the sample data covariance matrix, assuming the variations in the measured field are Gaussian-distributed, are given by the Wishart distribution (Wishart, 1928), which generalises the χ^2 -distribution;

$$p(\widehat{\mathbf{M}} | \mathbf{M}, \nu, \eta) = \left(\frac{\nu^{\nu\eta/2} |\mathbf{M}|^{-\nu/2} |\widehat{\mathbf{M}}|^{\gamma/2}}{2^{\nu\eta/2} \Gamma_{\eta}[\nu/2]} \right) e^{-\frac{\nu}{2} \text{Tr} \widehat{\mathbf{M}} \mathbf{M}^{-1}}, \quad (14)$$

where $|\mathbf{M}| = \det \mathbf{M}$ is the determinant of \mathbf{M} ,

$$\Gamma_{\eta} \left(\frac{\nu}{2} \right) = \pi^{\eta(\eta-1)/4} \prod_{s=1}^{\mu} \Gamma \left[\frac{\nu}{2} + \frac{1-s}{2} \right] \quad (15)$$

is the multivariate Gamma function (see Appendix B1 for a definition), $\eta = N_D$ is the size of the data-vector, \mathbf{M} and $\widehat{\mathbf{M}}$ are $\eta \times \eta$ matrices, ν is again the number of degrees of freedom of \mathbf{M} , and $\gamma = \nu - \eta - 1$. We require that $\nu > \eta$, to ensure the estimated data covariance matrix is positive definite.

For a single data point, where $\eta = 1$, the Wishart distribution is the reduced- χ^2 distribution,

$$p(y|\nu) = \left(\frac{\nu}{2} \right)^{\nu/2} \frac{y^{\nu/2-1}}{\Gamma[\nu/2]} e^{-\nu y/2} \quad (16)$$

where $y = \widehat{M}_{11}/M_{11} = \chi^2/\nu$, with mean $\langle y \rangle = 1$ and variance $\sigma^2(y) = 2/\nu$.

The mean of the general Wishart distribution is

$$\langle \widehat{\mathbf{M}} \rangle = \mathbf{M}, \quad (17)$$

showing it is indeed an unbiased estimate of the covariance matrix, while the covariance of $\widehat{\mathbf{M}}$ is

$$\langle \Delta \widehat{M}_{ij} \Delta \widehat{M}_{mn} \rangle = \frac{1}{\nu} (M_{im} M_{jn} + M_{in} M_{jm}). \quad (18)$$

This result can also be derived from the Gaussian four-point function or directly from the Wishart distribution (see Appendix B2 where we calculate the characteristic function for the Wishart).

3.3 The precision matrix and Inverse-Wishart

The simplest estimator for the precision matrix is

$$\widehat{\Psi} = \nu \left[\sum_{\alpha} \Delta \mathbf{D}_{\alpha} \Delta \mathbf{D}_{\alpha}^t \right]^{-1}, \quad (19)$$

where ν is the number of degrees-of-freedom. This estimator follows an inverse, or inverted, Wishart distribution (see e.g., Press, 1982),

$$p(\widehat{\Psi} | \Psi, \nu, \eta) = \left(\frac{\nu^{\nu\eta/2} |\widehat{\Psi}|^{-\beta/2} |\Psi|^{\nu/2}}{2^{\nu\eta/2} \Gamma_{\eta}[\nu/2]} \right) e^{-\frac{\nu}{2} \text{Tr} \widehat{\Psi}^{-1} \Psi}, \quad (20)$$

where $\beta = \nu + \eta + 1$, and $\eta = N_D$ is the size of the data-vector. We derive the Inverse-Wishart distribution in Appendix B3.

For a single data point, $\eta = 1$, the Inverse-Wishart reduces to the inverse- χ^2 distribution,

$$p(x|\nu) = \left(\frac{\nu}{2} \right)^{\nu/2} \frac{x^{-\nu/2-1}}{\Gamma[\nu/2]} e^{-\nu/2x}, \quad (21)$$

where $x = M_{11}/\widehat{M}_{11} = \nu/\chi^2$. The mean of this distribution is

$$\langle x \rangle = \frac{\nu}{\nu - 2}, \quad \nu > 2 \quad (22)$$

and its variance is given by

$$\sigma^2(x) = \frac{2\nu^2}{(\nu - 2)^2(\nu - 4)}, \quad \nu > 4. \quad (23)$$

Immediately we see that the inverse distribution has different properties to χ^2 . Not only is the mean of the inverse-distribution biased high, $\langle x \rangle > 1$, but both the mean and variance can diverge.

We can understand the behaviour of the inverse- χ^2 by considering the underlying Gaussian field. If $\Delta \mathbf{D}_{\alpha}$ is the

one data point, sampled N_S times, the sample variance is $\hat{M}_{11} = \sum_{\alpha} \Delta D_{\alpha}^2 / \nu$. As the ΔD_{α} fields are Gaussian, they will fluctuate symmetrically around $\Delta D_{\alpha} = 0$. Squaring and summing will produce positive values with mean $\langle \Delta D_{\alpha}^2 \rangle = M_{11}$. However the sample variance will scatter around this, bounded from below by zero. When we invert the sample variance, some of the values which are close to zero will become arbitrarily large. As there are no compensating small values, these large values will bias the mean of the precision matrix high, skewing the distribution.

The expectation value of the sampled precision matrix is biased and, assuming the mean is unknown, given by (Kaufman, 1967; see Press, 1982; Anderson, 2003; see also Hartlap et al., 2007, for a first application to cosmology)

$$\langle \hat{\Psi} \rangle = \frac{N_S - 1}{N_S - N_D - 2} \Psi. \quad (24)$$

If $N_S > N_D + 2$ is not satisfied then the values of $\hat{\mathbf{M}}$ are not positive-definite and its inverse is undefined. If we do satisfy this condition then the bias on the inverse can be corrected to yield an unbiased estimate of the precision matrix given by

$$\Psi_{\text{unbiased}} = \frac{N_S - N_D - 2}{N_S - 1} \hat{\Psi}. \quad (25)$$

In fact, this estimator for the precision matrix is the only unbiased estimator.

The covariance of the sample precision matrix, again assuming the mean is unknown, is (Kaufman, 1967; see also Press, 1982, Matsumoto, 2011)

$$\begin{aligned} \langle \Delta \hat{\Psi}_{ij} \Delta \hat{\Psi}_{mn} \rangle = \\ A \left[2\Psi_{ij}\Psi_{mn} + (N_S - N_D - 2)(\Psi_{im}\Psi_{jn} + \Psi_{in}\Psi_{jm}) \right], \end{aligned} \quad (26)$$

where

$$A = \frac{(N_S - 1)^2}{(N_S - N_D - 1)(N_S - N_D - 2)^2(N_S - N_D - 4)}. \quad (27)$$

The second term in equation (26) has the same form as the covariance for the data covariance matrix of Gaussian-distributed variables, and dominates when $N_S \gg N_D + 2$. This arises for large numbers of realisations as the Central-Limit Theorem will tend to make the Inverse-Wishart Gaussian distributed. The first term in equation (26) arises from the shift in the biased mean of the precision matrix.

The covariance matrix of the sample precision matrix is also biased high. If uncorrected, it leads to an overestimate of the uncertainty in the precision matrix and parameter errors. We can correct for the biases in the mean and covariance of the sample precision matrix, assuming the number of degrees-of-freedom is known and $N_S > N_D + 2$. The unbiased covariance of the precision matrix can be found by substituting the prefactor, A , in equation (26) by a corrected factor;

$$A_{\text{corr}} = \frac{1}{(N_S - N_D - 1)(N_S - N_D - 4)}. \quad (28)$$

The unbiased variance of the elements of the estimated precision matrix is

$$\sigma_{\text{corr}}^2[\hat{\Psi}_{ij}] = A_{\text{corr}} [(N_S - N_D)\Psi_{ij}^2 + (N_S - N_D - 2)\Psi_{ii}\Psi_{jj}], \quad (29)$$

where no summation over repeated indices is implied. A useful expression is the trace of this variance, which is given by

$$\text{Tr } \sigma_{\text{corr}}^2[\hat{\Psi}] = \frac{2}{(N_S - N_D - 4)} \sum_i \Psi_{ii}^2. \quad (30)$$

In the limit that $N_S \gg N_D + 4$ the uncertainty in the precision matrix falls off like the uncertainty on the data covariance matrix, $\sqrt{2/N_S}$. However, unlike the sample covariance matrix, the uncertainty on the sample precision matrix diverges when the number of realisations is close to the number of data points, while the sample covariance matrix is singular for fewer realisations. Hence, we find two closely related requirements. The number of independent realisations used to estimate the data covariance and precision matrix, N_S , must be larger than the total number of data-points, N_D , being measured, $N_S > N_D + 2$, to allow a correction for the bias in the estimated precision matrix, and $N_S > N_D + 4$ to avoid a divergence in the error on the precision matrix. Hence the bias and covariance of the precision matrix only depend on the number-of-degrees of freedom, $N_S - N_D$, and the model precision matrix, and are independent of the details of the experiment. This is a very simple, and powerful, result.

4 AN EXAMPLE FROM COSMOLOGY

4.1 Simulating a Weak Lensing Survey

As an example of the bias and variance of the sample precision matrix for a cosmological survey, we consider a simulation of a weak lensing shear survey. In Appendix A3 we describe the weak lensing fields and power spectra. Here we consider a single shear field, with no $B(\beta)$ -modes. We generated $N_S = 100$ samples of a 10×10 square degree weak lensing survey, with a Gaussian random shear field at a single redshift which we used to estimate the mean, covariance and precision matrix of the shear power. The surface density of galaxies is $\bar{n}_2 = 30$ per square degree, with a median redshift of $z_m = 0.9$. We did this for sample sizes from $N_S = 10$ to $N_S = 100$. We repeated this 100 times to generate independent groups of the N_S samples to estimate the mean and variance of the covariance and precision matrices.

This numerical experiment is non-trivial, as it has some more realistic assumptions compared to our analysis. In the simulations the underlying shear field is Gaussian-distributed, rather than the shear power itself. In practise this is what we expect for the large angular scale shear field, while on small scales we expect the shear field to be nonlinear and non-Gaussian. This will test if, on large scales, the assumption of Gaussian-distributed power is justified.

The Gaussianity of the shear field for a single redshift ensures that the shear power covariance matrix is diagonal, $M_{ij} = M_{\ell\ell'} = M_{\ell}\delta_{\ell\ell'}^K$. The covariance matrix of the sample data covariance matrix is

$$\langle \Delta \hat{M}_{\ell} \Delta \hat{M}_{\ell'} \rangle = \frac{2}{N_S - 1} M_{\ell}^2 \delta_{\ell\ell'}^K, \quad (31)$$

while the covariance matrix of the unbiased precision matrix is

$$\langle \Delta \hat{\Psi}_{\ell} \Delta \hat{\Psi}_{\ell'} \rangle = 2A_{\text{corr}} [(N_S - N_D - 2)\Psi_{\ell}^2 \delta_{\ell\ell'}^K + \Psi_{\ell}\Psi_{\ell'}]. \quad (32)$$

The Gaussian part is diagonal, as expected, however the shifted term introduces off-diagonal terms. The diagonal of the precision matrix now propagates into every term in the precision covariance matrix, with the ratio of diagonal (Gaussian) to non-diagonal (shift-term) terms scaling as $N_S - N_D - 1$. The variance of the bias-corrected precision matrix is

$$\sigma_{\text{corr}}^2[\hat{\Psi}_\ell] = \left(\frac{2}{N_S - N_D - 4} \right) \Psi_\ell^2. \quad (33)$$

Taking the sum of this, or the trace of the covariance matrix, we recover equation (30). Gaussian-distributed weak lensing convergence power spectra, on different angular scales and between redshift-bins, has the covariance

$$\langle |\Delta \hat{C}^{\kappa\kappa}(\ell, z, z')|^2 \rangle = \frac{1}{f_{\text{sky}}(2\ell + 1)} [|C^{\kappa\kappa}(\ell, z, z')|^2 + C^{\kappa\kappa}(\ell, z, z) C^{\kappa\kappa}(\ell, z', z')], \quad (34)$$

where f_{sky} is the fraction of the sky covered by the survey. If the shear power is binned into logarithmic passbands we divide by $\ell \Delta \ln \ell = \Delta \ell$, the number of ℓ -modes in each passband.

4.1.1 Whitening the covariance matrix

In cosmological surveys the data values can span several orders of magnitude, and so the conditional number of the corresponding covariance matrix is very large. This can cause numerical instabilities in the inversion to estimate the precision matrix, as well as the failure of some non-standard estimators (see Section 7.2). If a good model of the covariance elements is known, one can whiten the covariance matrix (see e.g., Bond, et al., 1998) by rendering its diagonal elements close to unity. Denoting the model data covariance elements by M_{ij}^{mod} and introducing a transformation matrix, $\mathcal{T}_{ij} = (1/\sqrt{M_{ii}^{\text{mod}}})\delta_{ij}^K$, as the inverse of the square-root of the diagonal elements of the model covariance, we define a new, whitened, data covariance matrix by $\hat{\mathbf{M}}_W = \mathcal{T} \hat{\mathbf{M}} \mathcal{T}$. The whitened matrix, \mathbf{M}_W , has a conditional number close to unity and is readily inverted, which we carry out using Singular Value Decomposition (SVD). The precision matrix is then obtained from the inverse of the whitened data covariance via $\Psi = \mathcal{T}^{-1} \mathbf{M}_W^{-1} \mathcal{T}^{-1}$. Here we whiten all data covariance matrices before correcting for whitening in the precision matrix.

4.1.2 Numerical Results

Figure 1 shows the mean and scatter (diagonal of the covariance) in the estimated shear power spectrum from our $N_S = 100$ simulated weak lensing survey, compared to the Λ CDM input model power. We have chosen $N_D = 36$ data points on the power spectrum as a compromise between oversampling the power spectrum with correlated data points, and under-sampling and missing some of its features which will contain parameter information.

Figure 2 shows the fractional bias in the trace of the mean of the sample precision matrix, which we define as

$$B = \frac{\sum_\ell \langle \hat{\Psi}_\ell \rangle - \sum_\ell \Psi_\ell}{\sum_\ell \Psi_\ell} = \frac{N_S - 1}{N_S - N_D - 2}, \quad (35)$$

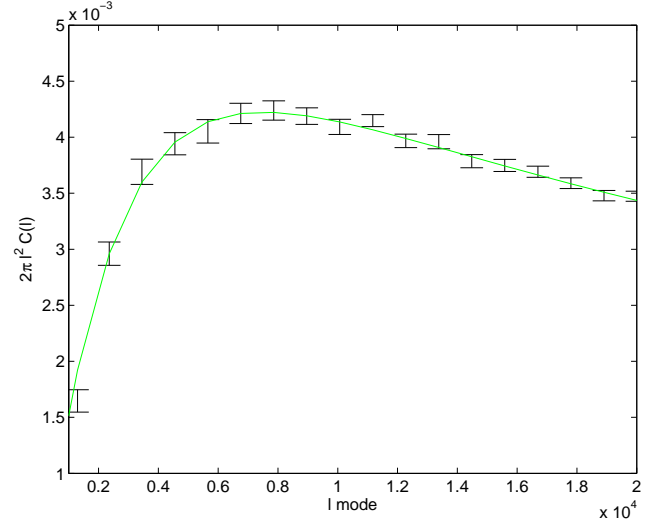


Figure 1. Simulated weak lensing shear auto-power spectrum, $C^{\kappa\kappa}(\ell, z, z)$ from 100 simulated 100 square degree surveys, for a single median redshift of $z = 0.9$. The solid line is the input power spectrum, while the data points are the mean estimated power spectrum, and the error bars are estimated from the 100 samples. We repeated this set of simulations 100 times to estimate the covariance of the data covariance and precision matrices.

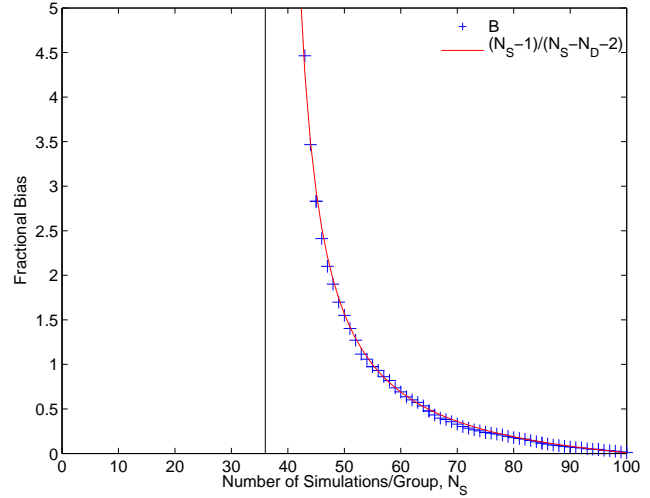


Figure 2. The fractional bias in the precision matrix of shear power spectra from N_S simulated and independent realisations of a 10^2 square degree weak lensing survey, with $N_D = 36$ power spectra data-points. The statistical properties of the precision matrix are generated from groups of 100 simulated surveys. The solid red line is the predicted scaling, while the vertical black line is the expected divergence for $N_S = N_D + 2$. Blue crosses are the estimated bias from the simulations, using equation (19), which closely follow the prediction. We have suppressed error bars on points, which are correlated.

where $\hat{\Psi}$ is estimated from equation (19), for $N_D = 36$ shear power spectra passbands as a function of number of realisations, N_S , in each group. The N_S realisations are cumulative, so each point is correlated with points on the left.

The numerical model, including the predicted divergence at $N_S = N_D + 4$ (solid vertical line), agrees ex-

tremely well with the prediction, as has previously been shown by Hartlap et al. (2007) in a similar cosmological context. There, they showed the bias followed the expected behaviour using a simulated weak lensing survey based on ray-tracing through many lines-of-sight of the Millennium N-body simulation, when the effects of non-linear clustering are included. The agreement between the simulations and prediction implies that the estimated precision matrix can indeed be debiased with the correction given by equation (25).

Figure 3 shows the measured fractional scaling of the trace of the covariance of the sample precision (crosses), defined as

$$E_\Psi = \sqrt{\frac{\sum_\ell \sigma^2(\hat{\Psi}_\ell)}{\sum_\ell \Psi_\ell^2}} = \sqrt{\frac{2}{N_S - N_D - 2}}, \quad (36)$$

and the trace of the fractional error on the data covariance matrix (stars), defined as

$$E_M = \frac{\sum_\ell \sigma(\hat{M}_\ell)}{\sum_\ell M_\ell} = \sqrt{\frac{2}{N_S - 1}}, \quad (37)$$

as a function of number of realisations, N_S . We have plotted the Wishart prediction (solid and dotted lines) for both statistics. The data covariance can be estimated for $N_S > 1$ data points, and its variance is stable below the $N_S = N_D$ line. However, the variance of the precision matrix estimate diverges when we reach the number of data points, as predicted by the Inverse-Wishart distribution. Again the sample of N_S realisations is cumulative and so each point is correlated to the points on its left. Again there is a very good agreement between the predicted and measured scaling of the variance of the data covariance matrix. The agreement between the predicted and numerical scaling of the variance of the precision matrix is also good, but there is some scatter and slight deviation which we attribute to the accuracy of the inversion of the sample data covariance.

4.2 The size of future surveys

As we have seen the main driver for the number of simulations comes from the precision matrix, whose accuracy is driven by the number of data points in our sample, $N_S > N_D + 4$. Here we discuss typical values which will be encountered. The issue of data compression will be discussed in Section 7.3.

4.2.1 Pixelised or discrete data-sets

In the case of pixelised data (Cosmic Microwave Background or Weak Lensing), or data where the individual data is sampled (Galaxy Redshift Surveys or Weak Lensing again) the number of data points can rise quite rapidly. If there is no analytic model for the pixel or data covariance matrix, the cost of simulations can be prohibitively high for $N_S > N_D$. We are then forced into some form of data compression, such as, in Cosmology, the estimation of two-point, or a number of n-point, power spectra or correlations.

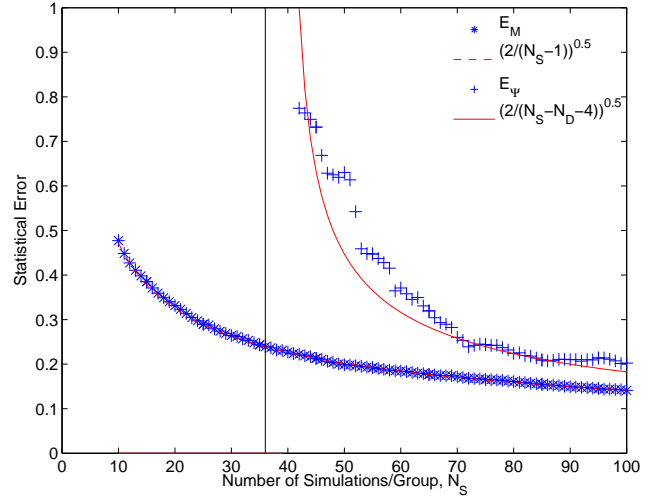


Figure 3. The error in the estimated precision and data covariance matrix from N_S realisations of the Weak Lensing power spectrum, with $N_D = 36$ data-points, generated from groups of 100 10^2 square degree simulated surveys, as a function of N_S . The vertical black line is the number of data points. The blue stars are the statistical errors on the unbiased data covariance matrix, compared to the predicted scaling (dotted line). Blue crosses are the statistical errors on the unbiased estimator of the precision matrix, equation (25), compared to the predicted scaling. Again the agreement is good, with differences due to numerical effects. We have suppressed error bars on points, which are again correlated.

4.2.2 Power spectrum analysis

In the case of a galaxy redshift survey, we imagine typically $N_k = 50$ data points sampling the galaxy spectra. To sample the anisotropic distortion in redshift-space we would need N_k^2 data points. If this was repeated in $N_b = 10$ redshift bins, we would have a total of $N_D = N_k^2 N_b = 2.5 \times 10^4$ data points in total.

For a tomographic, weak lensing power spectrum analysis measuring N_{spec} power spectra, over N_b -redshift bins, the total number of auto- and cross-spectra for spin-2 fields (including B -modes) is $N_b(2N_b + 1)$. If each power spectrum has N_ℓ passbands, the total number of data-points is $N_D = N_\ell N_b(2N_b + 1)$. If we again assume $N_b = 10$, we can measure 210 different power spectra. With $N_\ell = 50$ passbands per spectra per redshift, we have $N_D = 1.05 \times 10^4$ passbands, and the number of independent realisations we require is

$$N_S > 1.05 \times 10^4 \left(\frac{N_\ell}{50}\right) \left(\frac{N_b}{10}\right)^2. \quad (38)$$

If we add the lensing magnification power, $C^{\mu\mu}(\ell, z, z')$, to this, the estimated number of independent realisations increases by a significant factor.

If we combined cosmological probes we can estimate the size of a combined data-vector. If we assume $N_b = 10$ redshift bins, we have a total of $N_f = 43$ fields and

$$N_{\text{spec}} = \frac{1}{2} N_f (N_f + 1), \quad (39)$$

auto- and cross-spectra. For our example we then have

$N_{\text{spec}} = 946$ spectra. Assuming further that each spectrum has $N_\ell = 50$ passbands, we have

$$N_D = N_\ell N_{\text{spec}}, \quad (40)$$

or over 5×10^6 data-points in our data-vector. These raw numbers clearly represent a significant challenge for generating realisations. In Section 7 we discuss ways in which to avoid the Wishart bound and the need to generate such large numbers of simulated surveys.

5 PARAMETER COVARIANCE MATRIX

5.1 Covariance of the Fisher Matrix

Having found the statistical properties of the sample precision matrix we now turn to our main goal, to understand how the accuracy of the precision matrix propagates into maximum likelihood parameter estimation. To do this, we use the Fisher matrix formalism (e.g., Tegmark, Taylor & Heavens, 1997) to see how inaccuracies in the precision matrix leads to inaccuracy in the Fisher matrix and leads to inaccuracy in the parameter covariance matrix.

The log-likelihood, $\mathcal{L} \equiv -2 \log L$, can be expanded to second-order around its peak in parameter-space where the expectation value of the gradient of the log-likelihood is $\langle \partial_\alpha \mathcal{L} \rangle = 0$, while the expectation value of the curvature yields the Fisher Matrix,

$$\langle \partial_\alpha \partial_\beta \mathcal{L} \rangle = 2\mathcal{F}_{\alpha\beta}. \quad (41)$$

The derivatives of the mean are taken with respect to the parameters. For Gaussian-distributed data, this is given by (Tegmark, Taylor & Heavens, 1997)

$$\mathcal{F}_{\alpha\beta} = \frac{1}{2}(\partial_\alpha \boldsymbol{\mu} \partial_\beta \boldsymbol{\mu}^t + \partial_\alpha \boldsymbol{\mu}^t \partial_\beta \boldsymbol{\mu}) \boldsymbol{\Psi}. \quad (42)$$

With the Gaussian approximation, the likelihood surface of the parameter-space is specified completely by the parameter covariance matrix, Φ , given by the inverse of the Fisher matrix,

$$\Phi_{\alpha\beta} = \langle \Delta\theta_\alpha \Delta\theta_\beta \rangle = \mathcal{F}_{\alpha\beta}^{-1}. \quad (43)$$

If the data is again assumed Gaussian-distributed, the uncertainty on the precision matrix propagates into the Fisher matrix by

$$\Delta\mathcal{F}_{\alpha\beta} = \frac{1}{2}(\partial_\alpha \boldsymbol{\mu} \partial_\beta \boldsymbol{\mu}^t + \partial_\alpha \boldsymbol{\mu}^t \partial_\beta \boldsymbol{\mu}) \Delta\boldsymbol{\Psi}, \quad (44)$$

where $\Delta\boldsymbol{\Psi}$ is a random variation in the precision matrix. The covariance between terms in the Fisher matrices is given by

$$\begin{aligned} \langle \Delta\mathcal{F}_{\alpha\beta} \Delta\mathcal{F}_{\mu\nu} \rangle &= \frac{1}{4}(\partial_\alpha \boldsymbol{\mu} \partial_\beta \boldsymbol{\mu}^t + \partial_\alpha \boldsymbol{\mu}^t \partial_\beta \boldsymbol{\mu}) \langle \Delta\boldsymbol{\Psi} \Delta\boldsymbol{\Psi} \rangle \\ &\quad \times (\partial_\mu \boldsymbol{\mu}^t \partial_\nu \boldsymbol{\mu} + \partial_\mu \boldsymbol{\mu} \partial_\nu \boldsymbol{\mu}^t). \end{aligned} \quad (45)$$

We shall assume that the uncertainty in the mean of the data, $\boldsymbol{\mu}$, is negligible. Substituting equation (26) in for the covariance of $\boldsymbol{\Psi}$ we find the unbiased covariance of the Fisher matrix is

$$\begin{aligned} \langle \Delta\mathcal{F}_{\alpha\beta} \Delta\mathcal{F}_{\mu\nu} \rangle &= \\ A_{\text{corr}} [(N_S - N_D - 2) (\mathcal{F}_{\alpha\mu} \mathcal{F}_{\beta\nu} + \mathcal{F}_{\alpha\nu} \mathcal{F}_{\beta\mu}) + 2\mathcal{F}_{\alpha\beta} \mathcal{F}_{\mu\nu}] \end{aligned} \quad (46)$$

valid for Gaussian-distributed data.

5.2 Covariance of the parameter covariance

The parameter covariance matrix is the inverse of the Fisher matrix and so the uncertainty in the parameter covariance matrix is, to first-order,

$$\Delta\Phi_{\alpha\beta} = -\mathcal{F}_{\alpha\gamma}^{-1} \Delta\mathcal{F}_{\gamma\delta} \mathcal{F}_{\delta\beta}^{-1}, \quad (47)$$

where we assume summation over repeated indices. The covariance of the parameter covariance matrix is

$$\langle \Delta\Phi_{\alpha\beta} \Delta\Phi_{\mu\nu} \rangle = \Phi_{\alpha\delta} \Phi_{\eta\beta} \langle \Delta\mathcal{F}_{\delta\eta} \Delta\mathcal{F}_{\gamma\epsilon} \rangle \Phi_{\mu\gamma} \Phi_{\epsilon\nu}. \quad (48)$$

Substituting equation (46) for the covariance of the Fisher matrix, we find the covariance of the parameter covariance matrix is

$$\begin{aligned} \langle \Delta\Phi_{\alpha\beta} \Delta\Phi_{\mu\nu} \rangle &= A_{\text{corr}} [(N_S - N_D - 2) (\Phi_{\alpha\mu} \Phi_{\nu\beta} + \Phi_{\alpha\nu} \Phi_{\beta\mu}) \\ &\quad + 2\Phi_{\alpha\beta} \Phi_{\nu\mu}]. \end{aligned} \quad (49)$$

This is a central result of this paper. From this we see that the Inverse-Wishart covariance propagates through to the covariance of the parameter covariance matrix, with the Gaussian and shift terms. Again this diverges if $N_S \leq N_D + 4$. The components of this matrix can be written as

$$\langle |\Delta\Phi_{\alpha\alpha}|^2 \rangle = \frac{2}{N_S - N_D - 4} |\Phi_{\alpha\alpha}|^2, \quad (50)$$

$$\begin{aligned} \langle |\Delta\Phi_{\alpha\beta}|^2 \rangle &= A_{\text{corr}} [(N_S - N_D) r_{\alpha\beta}^2 \\ &\quad + (N_S - N_D - 2)] \Phi_{\alpha\alpha} \Phi_{\beta\beta}, \end{aligned} \quad (51)$$

$$\langle \Delta\Phi_{\alpha\alpha} \Delta\Phi_{\beta\beta} \rangle = 2A_{\text{corr}} [1 + (N_S - N_D - 2) r_{\alpha\beta}^2] \Phi_{\alpha\alpha} \Phi_{\beta\beta}, \quad (52)$$

$$\langle \Delta\Phi_{\alpha\beta} \Delta\Phi_{\beta\beta} \rangle = \frac{2}{N_S - N_D - 4} \Phi_{\alpha\beta} \Phi_{\beta\beta}, \quad (53)$$

where we have defined the parameter correlation coefficient,

$$r_{\alpha\beta} = \frac{\Phi_{\alpha\beta}}{\sqrt{\Phi_{\alpha\alpha} \Phi_{\beta\beta}}}. \quad (54)$$

From this result we see that the main factors which affect the covariance of the parameter covariances are the difference between the number of realisations of the survey and the size of the dataset, $N_S - N_D$, the degrees-of-freedom of the precision matrix, and the parameter correlation coefficient, $r_{\alpha\beta}$. This now provides us with a way to determine the accuracy with which we can estimate the distribution of parameter values in parameter space.

5.3 Accuracy of parameter errors

We can use this result to demonstrate how the accuracy of the errors on a parameter can be translated into the number of degrees-of-freedom in the data covariance, or equivalently the number independent realisations needed to estimate the precision matrix, and on the accuracy on the precision and data covariance matrices. For a single cosmological parameter (marginalised over all other parameters), the error on the parameter variance is given by equation (50). Comparing this to equation (33), and assuming that the data covariance is diagonal, we see that the fractional accuracy of the parameter variance is equal to the fractional accuracy of the precision matrix.

Defining ε as the fractional accuracy of the parameter variance (equation 50),

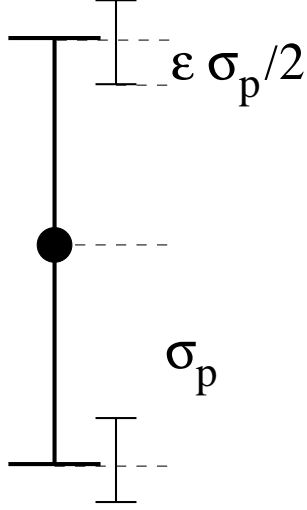


Figure 4. Sketch of the error on a parameter, p , given by σ_p , and the error on the error bar, $\epsilon\sigma_p/2$, where ϵ is the fractional variance on the precision matrix.

$$\epsilon = \frac{\sigma[\Phi_{\alpha\alpha}]}{|\Phi_{\alpha\alpha}|} = \sqrt{\frac{2}{N_S - N_D - 4}}. \quad (55)$$

we can consider an arbitrary parameter, p , with expected value p_0 . A measurement of p will yield an uncertainty, σ_p , while the error on that uncertainty will be $\epsilon\sigma_p/2$. We can write this as

$$p = p_0 \pm \sigma_p \left(1 \pm \frac{1}{2}\epsilon\right). \quad (56)$$

Figure 4 shows a sketch of this, illustrating for a single parameter the error bar and error on the parameter error.

The fractional error on the diagonal elements of the precision matrix is then given by (from equation 29, see also equation 33),

$$\frac{\sigma[\Psi_{ii}]}{|\Psi_{ii}|} = \epsilon. \quad (57)$$

Independent of the details of the survey, we find the required number of independent realisations of the survey for a given parameter accuracy and number of data points is

$$N_S > \frac{2}{\epsilon^2} + (N_D + 4). \quad (58)$$

The first term here is **the** usual root- N_S scaling for independent samples, and sets a lower limit on the number of independent realisation required to reach a given accuracy. The second term arises from the Inverse-Wishart variance, where the number of independent realisations needed to **reach** a given accuracy scales as the number of data points. Figure 5 shows the scaling of the number of independent samples, N_S , with the size of the data-set, N_D . The value of N_S in the limit $N_D \rightarrow 0$, is set by the desired accuracy of the parameter variance.

The fractional error on the data covariance matrix, for a given parameter error accuracy and number of data points, is

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} < \sqrt{\frac{2\epsilon^2}{2 + \epsilon^2(N_D + 4)}}. \quad (59)$$

This has two regimes. When $\epsilon^2 \ll 2/(N_D + 4)$, i.e. for small

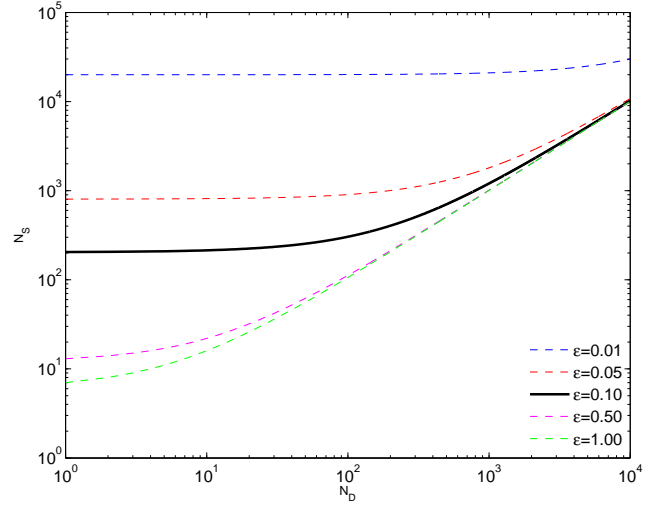


Figure 5. Scaling of the number of independent realisations of the survey, N_S , as a function of the size of the data set, N_D , for different fractional accuracies on the variance of the parameter variance, ϵ .

data-sets compared to the required accuracy, the fractional error on the data covariance scales as

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} = \epsilon, \quad (60)$$

the same as for the precision matrix and the variance on the parameter variance, while for $\epsilon^2 \gg 2/(N_D + 4)$, when the data-set is large, the error on the data covariance scales as

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} < \sqrt{\frac{2}{N_D + 4}} \ll \epsilon. \quad (61)$$

For large-data sets, this scaling puts the strongest constraints on the accuracy of the data covariance matrix, which can be much higher than the accuracies of the precision and parameter covariance matrices.

5.4 Constraining the parameter error

While the accuracy of the parameter error is set by the number of independent realisations of the survey, N_S , and the data size, N_D , it is useful to consider what typical accuracies any analysis should achieve, independent of the details of the particular survey. A reasonable accuracy for a parameter error is 5%, since a much higher accuracy will put strong requirements **on** the number of realisations, while lower accuracy will compromise the measurement error. This requires that the marginalised parameter variance should be accurate to 10% and that the precision matrix, Ψ , is accurate to 10%, or $\epsilon = 0.1$. From equation (58), this requires $N_S > 200 + N_D + 4$ independent realisations of the survey to reach this accuracy.

For a small data set, with $N_D \ll 100$, a minimum of 204 independent realisations of the survey yields a 5% error on the parameter error. This implies that the data covariance matrix is accurate to 5%. When the data-set becomes $N_D \gg 100$, we require $N_S > N_D + 4$ independent realisations, and the fractional accuracy of the data covariance matrix scales as

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} \approx \sqrt{\frac{2}{N_D}} \ll 5\%. \quad (62)$$

In particular, for forthcoming tomographic cosmological surveys with 10 redshift bins, we can expect some 10^4 power spectra data points requiring at least this number of independent realisations. This will increase the accuracy on the data covariance matrix to 1.4%. For combined data-sets the number of data points can rise to $N_D \approx 10^6$, and hence require an accuracy of 0.14% on the data covariance matrix. Achieving these accuracies will be challenging.

6 FIGURES-OF-MERIT

6.1 Uncertainty of the Figure-of-Merit

In addition to marginalised parameter errors, it is useful to know how the uncertainty in the precision matrix affects the FoM. The FoM is the inverse of the enclosed area within a certain likelihood contour, and is frequently used as a target statistic to optimise cosmological surveys. Our aim here is to understand how inaccuracies in the precision matrix propagate through the parameter estimation into a FoM, and how fixing the required FoM can be used to put constraints on the accuracy of the precision matrix and data covariance matrix.

The dark energy Figure-of-Merit (DE FoM), $\Xi_{w_0 w_a}$, is defined as the inverse of the area of the 68% error-ellipse for a two-parameter dark energy model (e.g., Albrecht et al., 2006),

$$\Xi_{w_0 w_a} = \frac{1}{\sqrt{\Phi_{w_0 w_0} \Phi_{w_a w_a} - \Phi_{w_0 w_a}^2}}, \quad (63)$$

where w_0 and w_a parameterise the dark energy equation of state, $w(a) = \rho_{de}(a)/P_{de}(a) = w_0 + w_a(1-a)$ (Chevallier & Polarski, 2001; Linder, 2003), where $a(t)$ is the cosmological scale factor, and ρ_{de} and P_{de} are the energy-density and pressure of the dark energy. We define a general FoM matrix for any two parameters as

$$\Xi_{\alpha\beta} = \frac{1}{\sqrt{\Phi_{\alpha\alpha}\Phi_{\beta\beta} - \Phi_{\alpha\beta}^2}} = \frac{1}{\sqrt{(1-r_{\alpha\beta}^2)\Phi_{\alpha\alpha}\Phi_{\beta\beta}}}, \quad (64)$$

where no summation over the repeated indices α and β is implied. In the second expression we have used the parameter correlation coefficient, $r_{\alpha\beta}$.

It is useful to consider the inverse of the elements of the FoM, the area of each ellipse in the parameter space

$$A_{\alpha\beta} = \frac{1}{\Xi_{\alpha\beta}}. \quad (65)$$

The fractional change in $\Xi_{\alpha\beta}$ due to a change in the area is

$$\frac{\Delta\Xi_{\alpha\beta}}{\Xi_{\alpha\beta}} = -\frac{\Delta A_{\alpha\beta}}{A_{\alpha\beta}}. \quad (66)$$

Varying the parameter covariance matrix in the FoM, we find the fractional change in the area of the error ellipse is

$$\frac{\Delta A_{\alpha\beta}}{A_{\alpha\beta}} = \frac{1}{2(1-r_{\alpha\beta}^2)} \left(\frac{\Delta\Phi_{\alpha\alpha}}{\Phi_{\alpha\alpha}} + \frac{\Delta\Phi_{\beta\beta}}{\Phi_{\beta\beta}} - 2r_{\alpha\beta} \frac{\Delta\Phi_{\alpha\beta}}{\sqrt{\Phi_{\alpha\alpha}\Phi_{\beta\beta}}} \right). \quad (67)$$

Using the results of the covariance of the parameter covariance matrix, equations (50) to (53) for the unbiased precision matrix, we find the variance of each FoM is

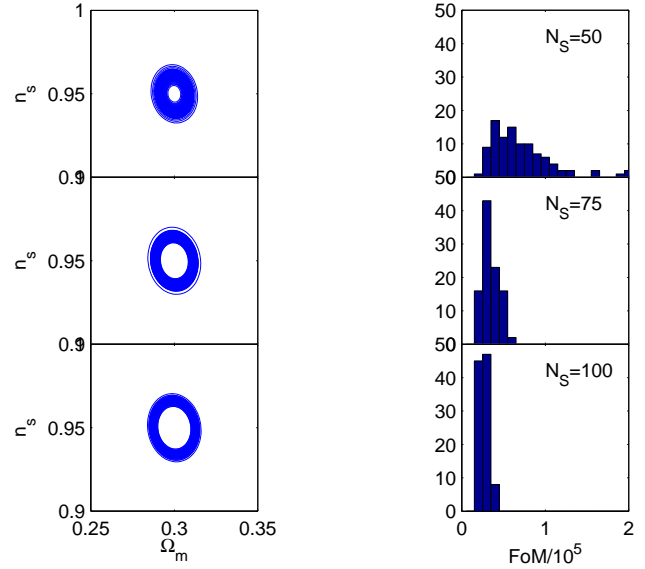


Figure 6. Numerical simulation of the parameter likelihood contours, estimated from a simulated weak lensing survey. The left-hand side (LHS) shows the 2-parameter 86.3% likelihood contour for the cosmological parameters, Ω_m , the density parameter of matter and the clustering spectral index n_s . The right-hand-side (LHS) shows the frequency distribution of the FoM (inverse area) for these parameters over 100 realisations. The top row is for a sample-size of $N_S = 50$ realisations, the middle row for $N_S = 80$, and the bottom row for $N_S = 100$. As the number of realisations increases the accuracy increases as predicted.

$$\sigma^2[\Xi_{\alpha\beta}] = \frac{(N_S - N_D)}{(N_S - N_D - 4)(N_S - N_D - 1)} |\Xi_{\alpha\beta}|^2. \quad (68)$$

Once again, the fractional variance of the FoM depends only on the difference between the number of independent realisations of the survey used to estimate \mathbf{M} and $\mathbf{\Psi}$, and the size of the data set.

In Figure 6 (LHS) we show the 2-parameter, marginalised likelihood surface in the $\Omega_m - n_s$ plane for our weak lensing simulations. Each ellipse is a 2-parameter, 68.3% likelihood contour for the group of simulations with N_S realisations. As the number of realisations increases from $N_S = 50$ to 100, we see the spread in areas decreases. To quantify this, in Figure 6 (RHS) we plot the frequency distribution of the FoM for this parameter plane.

Figure 7 shows the predicted uncertainty in the FoM, equation (68), compared to the variance of the FoM distributions shown in the RHS column in Figure 6, as a function of number of realisations, N_S . We see good agreement between our prediction and the error measured on the FoM from the weak lensing simulation.

6.2 Accuracy of the Figure-of-Merit

In the design of many cosmological surveys, the FoM is used as a target statistic to optimise the survey design, varying area, depth and number of photometric passbands to find the design which maximises the FoM. Having set this optimal FoM, we then want to keep biases and uncertainties down to a level which does not violate the expected FoM. Here we develop an approach which uses the required FoM as

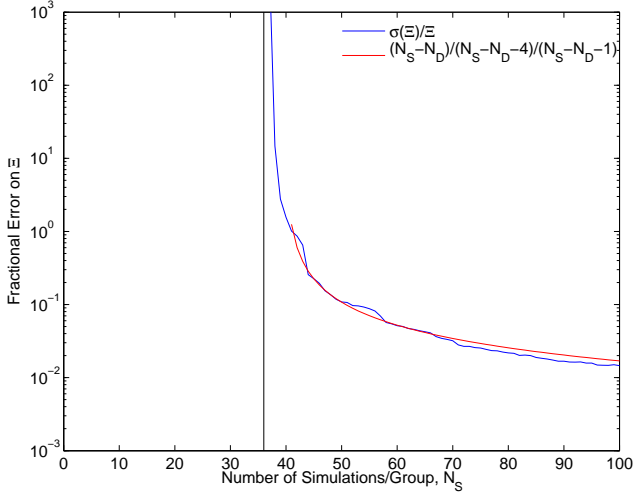


Figure 7. The scaling of the fractional error on the FoM, $\Xi_{\alpha\beta}$, as a function of the number of simulated realisations of a weak lensing survey, N_S , with $N_D = 36$ data points. The solid blue line is the scaling predicted from the Inverse Wishart distribution, while the red line is the scaling found from the simulated surveys.

a constraint to determine the number of survey realisations needed to do this. We then translate this into the accuracies of the precision and data covariance matrices.

In order to keep within a required FoM we set the constraint that the uncertainty in the likelihood area, when added in quadrature with the area, should not exceed some fiducial value, $A_{\alpha\beta}^0$,

$$A_{\alpha\beta}^2 + \sigma^2[A_{\alpha\beta}] \leq (A_{\alpha\beta}^0)^2. \quad (69)$$

Assuming the uncertainty in the area is small, and taking the expectation value, we can re-write this in terms of the FoM,

$$\Xi_{\alpha\beta} \left(1 - \frac{1}{2} \left(\frac{\sigma[\Xi_{\alpha\beta}]}{|\Xi_{\alpha\beta}|} \right)^2 \right) \geq \Xi_{\alpha\beta}^0. \quad (70)$$

The effect of a random change in the area of the error ellipse in parameter-space will, on average, reduce the FoM. Hence the actual FoM we need to measure, $\Xi_{\alpha\beta}$, to meet the required $\Xi_{\alpha\beta}^0$ is increased. We define the fractional error in the FoM as

$$\varepsilon_{\Xi} \equiv \frac{\sigma[\Xi_{\alpha\beta}]}{|\Xi_{\alpha\beta}|}. \quad (71)$$

In order to keep the fractional increase in the FoM below some value, ε_{Ξ}^2 , we can solve equation (68) for N_S and find that the number of independent realisations should be

$$N_S > N_D + \frac{5}{2} + \frac{1}{2\varepsilon_{\Xi}^2} \left(1 + \sqrt{(1 + 9\varepsilon_{\Xi}^2)(1 + \varepsilon_{\Xi}^2)} \right). \quad (72)$$

In the limit that $\varepsilon_{\Xi} \ll 1$ we find,

$$N_S > N_D + \frac{1}{\varepsilon_{\Xi}^2}. \quad (73)$$

If we want the fractional error on the FoM to be 10%, the number of realisation required (using equation 72) is

$$N_S > N_D + 125. \quad (74)$$

Again, we see that for small data-sets the number of realisations is fixed, this time at $N_S = 125$, while for large-data-sets the number of realisations again scales as the number of data points.

6.3 Accuracy of the precision matrix

To set a constraint on the accuracy of the precision matrix, for a given accuracy on the FoM, we again only consider the diagonal components of the precision matrix, where $\sigma[\Psi_{ii}] = \varepsilon|\Psi_{ii}|$ (see equation 57). Substituting the constraint from the FoM on the number of realisations, equation (72), we find

$$\frac{\sigma^2[\Psi_{ii}]}{|\Psi_{ii}|^2} = \frac{4\varepsilon_{\Xi}^2}{1 + \sqrt{(1 + 9\varepsilon_{\Xi}^2)(1 + \varepsilon_{\Xi}^2)} - 3\varepsilon_{\Xi}^2}, \quad (75)$$

which only depends on the accuracy of the FoM. In the high-accuracy regime, $\varepsilon_{\Xi} \ll 1$, this reduces to

$$\frac{\sigma[\Psi_{ii}]}{|\Psi_{ii}|} \approx \sqrt{2}\varepsilon_{\Xi}. \quad (76)$$

If we require for the FoM that $\varepsilon_{\Xi} = 0.1$ this implies that

$$\sigma[\Psi_{ii}] \approx 0.19|\Psi_{ii}|, \quad (77)$$

or an accuracy of 19% on the precision matrix.

6.4 Accuracy of the data covariance matrix

Given we know that

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} = \sqrt{\frac{2}{N_S}}, \quad (78)$$

and that the number of independent realisations required to reach a given accuracy of the FoM scales according to equation (72), we can write

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} = \frac{2\varepsilon_{\Xi}}{\sqrt{1 + (5 + 2N_D)\varepsilon_{\Xi}^2 + \sqrt{(1 + 9\varepsilon_{\Xi}^2)(1 + \varepsilon_{\Xi}^2)}}}. \quad (79)$$

In the high-accuracy regime, $\varepsilon_{\Xi}^2 \ll 1$, this reduces to

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} \approx \sqrt{\frac{2\varepsilon_{\Xi}^2}{1 + N_D\varepsilon_{\Xi}^2}}, \quad (80)$$

which for $N_D\varepsilon_{\Xi}^2 \ll 1$ reduces further to $\approx \sqrt{2}\varepsilon_{\Xi}$, scaling like the fractional accuracy of the precision matrix in the same regime, while for $N_D\varepsilon_{\Xi}^2 \gg 1$ the fractional error reduces to $\sqrt{2/N_D} \ll \sqrt{2}\varepsilon_{\Xi}$. Hence, for high-accuracy FoM's and large data-sets, the accuracy of the data covariance matrix is driven by the size of the data-set.

For our fiducial accuracy of $\varepsilon = 0.1$ we find

$$\frac{\sigma[M_{ii}]}{|M_{ii}|} \approx \sqrt{\frac{0.02}{1 + (N_D/100)}}, \quad (81)$$

where for $N_D \ll 100$, the error on the data covariance is $\sigma[M_{ii}] \approx 0.19|M_{ii}|$, the same accuracy as the precision matrix, while for $N_D \gg 100$ we find the $\sqrt{2/N_D}$ scaling.

7 BEYOND THE WISHART BOUND

The main conclusion of our analysis is that without an accurate model data covariance matrix, and if we sample the data

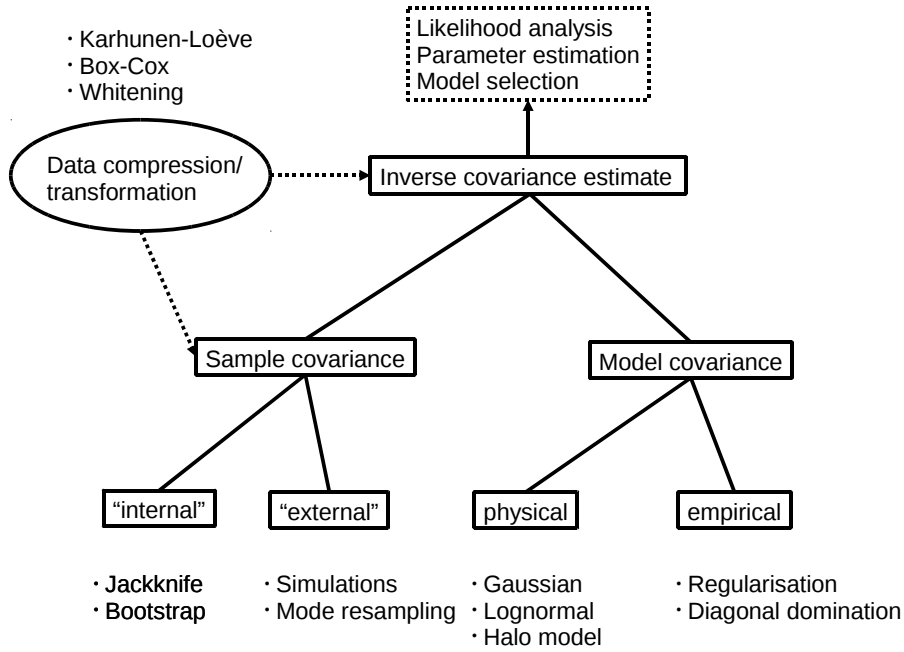


Figure 8. Different routes to determining inverse covariance matrices for use in likelihood analysis. The bottom level shows different ways to estimate sample and model covariance matrices, leading to the precision matrix and the likelihood function.

covariance, the number of independent realisations needs to be greater than the number of data-points we are analysing. For large data-sets, such as for the surveys now underway in Cosmology, this requires a prohibitively large number of simulations of the data. In this Section we discuss alternative routes to obtaining an accurate precision matrix which may help to meet, or avoid, the tight requirements set by simple estimation. Alternative approaches will also provide valuable consistency checks on some of our assumptions. We consider four methods: theoretical modelling of the data covariance, optimal estimators, data compression, and simulation and data resampling. Figure 8 shows a schematic view of these methods and their relationship. We shall not consider more radical alternatives, such as going directly to estimates of the likelihood function itself.

7.1 Theoretical modelling of the data covariance

The problems of noise, inherent to the sample covariance matrix, are avoided altogether if we can accurately model the data covariance matrix analytically. Modelling of the data covariance ranges from assuming the data is Gaussian distributed on large-scales (e.g. Kaiser, 1992; Knox, 1995) to assessing the impact of non-linear clustering, using perturbation theory and simulations, on the galaxy power spectrum covariances (Meiksin & White, 1999) and the weak lensing power spectra covariance (Scoccimarro, Zaldarriaga, Hui, 1999). In addition, the halo model has been used to estimate the covariance matrix for the matter and weak lensing power spectra (e.g., Cooray & Hu, 2001; Takada & Bridle, 2007; Takada & Jain 2009; Kayo et al., 2012), while estimation of the data covariance matrix for a lognormal field (e.g., Hilbert, et al., 2011) seems to reproduce the main features of the covariance structure for weak lensing correlations. Given

the level of difficulty in modelling the nonlinear regime one needs to test the range of validity of these model against simulations. For example, Takahashi et al. (2009) have tested nonlinear modelling of the galaxy clustering data covariance matrix on a suite of 5000 simulations, while Sato et al. (2009) have tested nonlinear estimates of the weak lensing covariance on simulations. Kiessling et al. (2011) have also studied modelling non-Gaussian covariances and parameter forecasting using weak lensing simulations. Hamilton & Rimes (2005, 2006) first pointed out the loss of information in the quasi-nonlinear regime in the matter power spectrum due to non-Gaussianity using simulations, while Hamilton, Rimes & Scoccimarro (2005) have discussed some of the issues with estimating the data covariance matrix from simulations. Theoretical models of the data covariance matrix have been applied to parameter estimation from data with the CMB (e.g. Verde, et al., 2003; Spergel, et al., 2003), galaxy redshift surveys (e.g. Ballinger, et al., 1995; Tadros, et al., 1999), and weak lensing (e.g., Brown, et al., 2003; Kitching, et al., 2007).

Even if modelling is not precise, one could develop analytic parameterised fitting functions to the data covariance, fitting the free parameters to simulations or the data. Importantly, these physically motivated models could allow one to incorporate the cosmology dependence of the covariance, which would require a substantial increase in the number of simulated realisations to cover parameter space. Sometimes, truncation or smoothing of the sample data covariance is used to suppress noise (e.g. Mandelbaum et al., 2012). However, such approaches alter the number of degree-of-freedom in the data covariance and so we would no longer know how to correct the precision matrix for bias.

If we assume for the moment that we can model the theoretical uncertainty on the data covariance matrix as

random, even for the next generation of surveys when we expect $N_D \approx 10^4$, the data covariance has to be known to a few percent accuracy, which will become a problem for theoretical computation. For surveys with 10^6 data points the accuracy of the data covariance has dropped to a fraction of a percent, putting high demands on its calculation (see Section 5.4).

7.2 Optimal precision estimators

If we do not have a reliable model or fitting function to the data covariance, we can still suppress the noise in the sample covariance by combining it with some simple model or prior knowledge. This is generally referred to as shrinkage estimation. One can define optimised covariance estimators in the sense that they yield smaller variance than equation (29) while keeping the bias small. We investigate the performance of three well-known cases.

7.2.1 Shrinkage: Stein precision estimator

Stein et al. (1972) have proposed the precision estimator

$$\hat{\Psi}_{\text{Stein}} = \frac{N_S - N_D - 2}{N_S - 1} \hat{\Psi} + \frac{N_D(N_D + 1) - 2}{(N_S - 1) \text{Tr} \hat{\mathbf{M}}} \mathbf{I}, \quad (82)$$

which is defined for $N_S > N_D + 2$. If $N_S \gg N_D$, the estimator reduces to the unbiased estimate, $\hat{\Psi}_{\text{unbiased}}$. If $N_S \sim N_D \gg 1$, the estimator returns $S^{-1} \mathbf{I}$, where

$$S = \frac{1}{N_D} \text{Tr} \hat{\mathbf{M}} \quad (83)$$

is the average of the diagonals of $\hat{\mathbf{M}}$. This is exact if the covariance is diagonal and homoscedastic. This estimator has smaller loss than any estimator that is proportional to $\hat{\Psi}$ for a ‘natural’ loss function (Stein, et al., 1972), a generalisation of least squares between the matrix elements of the estimator and the true precision matrix.

7.2.2 Shrinkage: Haff precision estimator

A second estimator, suggested by Haff (1974), is;

$$\hat{\Psi}_{\text{Haff}} = \frac{N_S - N_D - 2}{N_S - 1} \left((1 - \sqrt{U}) \hat{\Psi} + \sqrt{U} S^{-1} \mathbf{I} \right), \quad (84)$$

with

$$U = S^{-1} |\hat{\mathbf{M}}|^{1/N_D} \quad (85)$$

again defined for $N_S > N_D + 2$ only. The variable U measures disparity among the eigenvalues of $\hat{\mathbf{M}}$ and lies in the interval $0 \leq U \leq 1$, shifting from the unbiased sample estimator ($U = 0$) to the estimator $\propto S^{-1} \mathbf{I}$ ($U = 1$). This estimator has smaller loss than any estimator that is proportional to $\hat{\Psi}$ for a whole class of loss functions (Haff, 1974). However, this property is only guaranteed close to the divergent case, in our case for $N_S \leq N_D + 4$.

7.2.3 Target data covariance shrinkage

Shrinkage in its narrower sense refers to covariance estimates in which the balance between the sample covariance and the assumed model (the ‘target’, see Section 7.1) is estimated

from the data as well. The estimate for the precision matrix is given by the inverse of

$$\hat{\mathbf{M}}_{\text{shrink}} = \lambda \mathbf{T} + (1 - \lambda) \hat{\mathbf{M}}, \quad (86)$$

where \mathbf{T} is the theoretical target covariance matrix. This formalism has been applied to covariance estimation of galaxy clustering power spectra by Pope & Szapudi (2008), and to the CMB by Hamimeche & Lewis (2009). Ledoit & Wolf (2003) derived an analytic estimator for the shrinkage intensity λ , thereby greatly reducing the computational cost of this form of shrinkage estimation. It is given by (see also Schäfer & Strimmer, 2005)

$$\lambda = \frac{\sum_{ij}^{N_D} \text{Var}[\hat{M}_{ij}] - \text{Cov}[T_{ij}, \hat{M}_{ij}]}{\sum_{ij}^{N_D} \text{Var}[\hat{M}_{ij} - T_{ij}] + (\bar{\hat{M}}_{ij} - \bar{T}_{ij})^2}, \quad (87)$$

where the variances and covariances are computed from the N_S realisations, so then (see also Pope & Szapudi, 2008)

$$\text{Var}[\hat{M}_{ij}] = \frac{N_S^2}{(N_S - 1)^3} \sum_{\alpha=1}^{N_S} \left(W_{ij}^{(\alpha)} - \bar{W}_{ij} \right)^2, \quad (88)$$

where

$$W_{ij}^{(\alpha)} = \Delta D_{\alpha,i} \Delta D_{\alpha,j}. \quad (89)$$

As only one set of realisations is available in practice, the means in the denominator of Equation (87) have to be replaced with estimates of $\hat{\mathbf{M}}$ and \mathbf{T} . If the target matrix is noise-free, $\text{Cov}[T_{ij}, \hat{M}_{ij}] = 0$ and $\text{Var}[\hat{M}_{ij} - T_{ij}] = \text{Var}[\hat{M}_{ij}]$. If $\bar{\hat{M}}_{ij} - \bar{T}_{ij} = 0$, the target accurately describes the covariance in the data and $\lambda = 1$. Conversely, the shrinkage intensity tends to zero if the target attains a similar noise level as, and/or if the mean target deviates strongly from the mean of, the sample data covariance.

We consider two choices for our target matrix to test on our weak lensing simulations. The first is a theoretical estimate of the covariance matrix based on a Gaussian-distributed power spectrum (e.g., Kaiser, 1992),

$$\mathbf{T}_1 = \left(\frac{C_i^2}{f_{\text{sky}} \ell_i^2 \Delta \ln \ell} \right) \mathbf{I}, \quad (90)$$

where f_{sky} is the fraction of the sky covered by the survey, and we have assumed log-binning. Since this should correspond closely to the simulations we have generated we expect $\lambda \rightarrow 1$. The second is an empirical target matrix,

$$\mathbf{T}_2 = S \mathbf{I}, \quad (91)$$

which represents a minimum-knowledge approach.

7.2.4 Testing Shrinkage

Figures 9 and 10 show the fractional bias and error on the precision matrix for the Stein estimator (black points), the Haff estimator (green points), and the two target shrinkage methods, model (blue points) and mean (purple points), applied to our weak lensing simulations.

• **Stein estimator:** The Stein estimator is more biased than the simplest sample estimator for large numbers of realisation, even though it asymptotes to become the same estimator. Simulations would be required to calibrate this

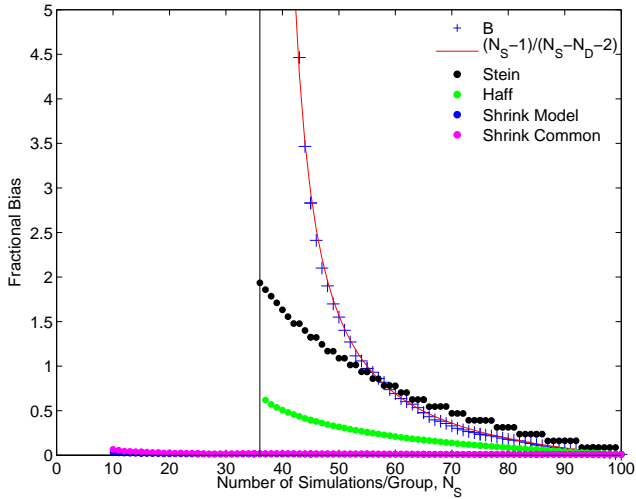


Figure 9. The bias in the estimated precision matrix from N_S realisations of the Weak Lensing power spectrum, with $N_D = 36$ data-points, generated from groups of 100×10^2 square degree simulated surveys, as a function of N_S . The blue circles are for direct inversion of the unbiased data covariance matrix, while green are for the Stein and Haff estimators.

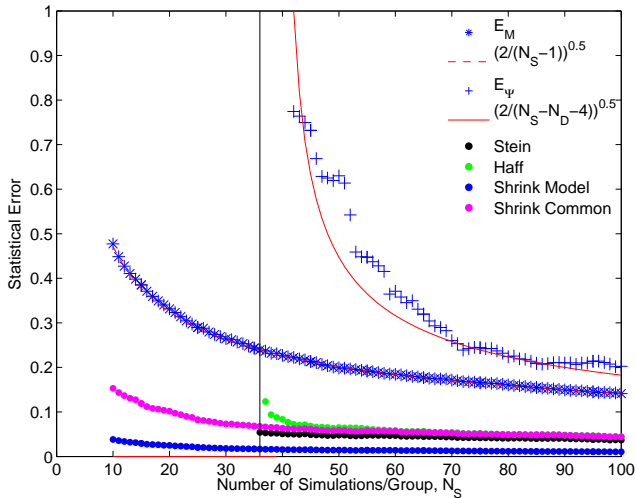


Figure 10. The error in the estimated precision matrix from N_S realisations of the Weak Lensing power spectrum, with $N_D = 36$ data-points, generated from groups of 100×10^2 square degree simulated surveys, as a function of N_S . The blue circles are for direct inversion of the unbiased data covariance matrix, while green are for the Stein and Haff estimators.

for any particular experiment. The variance of the Stein estimator is low, yielding an almost constant 5% error.

- **Haff estimator:** The Haff estimator is less biased than the Stein, and the sample estimator, but still shows significant bias. The variance of the Haff estimated precision matrix is again around 5%.

- **Target estimator:** The T_1 model target shrinkage estimator is essentially unbiased, as we expect for an accurate model, and works for $N_S < N_D + 4$. The T_1 model also does best at lowering the statistical to around 1% accuracy. The T_2 , empirical estimator yields a similarly unbiased precision

matrix, which is of interest. The error on the estimated precision matrix is slightly higher than the model target, at 5%, similar to that for the Stein and Haff estimators.

We conclude that the model estimator, T_1 , works impressively well when the model is a good approximation, minimising the bias in the precision matrix and reducing the error in the precision matrix to a few percent. For more realistic simulations, we would expect to take advantage of the more detailed theoretical models (see Section 7.1). Interestingly, the empirical T_2 target is similarly unbiased, with a 5% error, although we caution that our simulation covariance is also diagonal. Both of these estimators would satisfy our goal of an unbiased parameter errors with 5% error, requiring only $N_S \approx 30$ realisations compared to the $N_S > 240$ needed for the sample data covariance. We imagine this would work just as well for much larger data-sets. The Haff and Stein estimator have a bias similar to the sample estimator, but without the known correction factor. One would have to calibrate these with simulations. If this can be done, the error is sufficiently low to be useable.

7.3 Data compression

Independently of how the realisations to compute the sample covariance matrix are obtained, one can lower the Wishart bound by reducing the size of the data vector. Karhunen-Loève eigenvalue methods (Tegmark, Taylor & Heavens, 1997, and references therein) are widely used in astronomy to compress data by finding a smaller data vector that maximises the Fisher information (preserving parameter information), while simultaneously diagonalising the new data-vector covariance. Heavens, et al. (2000) introduced a linear implementation that is lossless for an arbitrary number of estimated parameters if the data covariance is not parameter-dependent, and otherwise still performs better than principal component analysis of the Fisher matrix. Data compression can also help stabilise the inversion of the data covariance, since the noisy modes which cause numerical instabilities are removed (e.g., Taylor, et al., 2001).

We can consider what optimal compression may achieve based on general considerations. The compression method of Heavens et al. (2000) compresses data to one element per non-degenerate cosmological parameter. Future surveys will constrain cosmologies with a large number of parameters together with a substantial list of calibration and nuisance parameters, so that one can expect several tens of parameters in total, corresponding to a compression to $N_D \sim 100$. Alternatively, one can argue from the number of characteristic features in the matter power spectrum, which features a number of transition scales beyond the overall amplitude and slope, and the clustering growth rate, leading to an estimate of several tens of parameters.

7.3.1 Compression in CMB analysis

Gupta & Heavens (2002) demonstrate that the some 400 temperature power spectrum modes can be compressed to the 10 – 20 cosmological parameters of interest. The computation of these modes requires a one-off $O(N_{\text{para}} N_D^3)$ calculation, compared to the $O(N_D^3)$ needed for a brute-force likelihood approach. However our aim here is to minimise

the uncertainty on the likelihood function, rather than to speed up the analysis.

7.3.2 Compression in weak lensing

In lensing the COSEBI two-point statistics are an efficient way of compressing angular weak lensing data. Asgari et al. (2012) found that of the order 10 modes capture the bulk of the cosmological information. Moreover, the B-mode is expected to have low signal-to-noise while cross-correlations between E- and B-mode should vanish altogether, so that these signals can be compressed efficiently into a small number of elements.

For radial lensing modes, Heavens, et al. (2003) applied Karhunen-Loève methods to a 3-D weak lensing analysis and found that only 4 (out of 100) radial modes contain significant information (see Hu, 1999, for similar conclusions on tomographic weak lensing data). Tomographic or 3-D weak lensing analyses will also have to simultaneously model intrinsic galaxy alignments, which are primarily separated from the lensing signal via their different redshift dependence. A maximum number of 10 radial elements in the data vector per angular frequency should be a conservative estimate. The main issue is then the independence of these modes, since we have of the order $\sim N_b^2$ cross-spectra. If the modes are independent, we only need around $N_D \sim 100$ spectra to consider, in agreement with the estimate for Karhunen-Loève methods. If they are not truly independent and we rely on angular compression, we may only compress weak lensing data by a factor of a few. This gives us a range of possible compression factor from 10 to a few.

If data compression can compress data so that $N_D \approx N_{\text{para}} \approx 100$, this would be very powerful and imply we only ever need around $N_S \approx 300$ realisations for any survey. However, this is probably over-optimistic, and if the real compression is a factor of 10, large data-sets of $N_D \approx 10^4 - 10^6$ will still be difficult to accurately analyse.

7.4 Resampling techniques

In previous Sections we have considered how we can reduce the need to generate large numbers of realisations of our surveys to ensure the accuracy of parameter errors. Here we discuss how we could generate these realisations. In general we can divide this into external realisations, usually from simulations, or internal realisations from e.g., Jackknife or Bootstrap resampling of the data itself.

7.4.1 Simulation Mode Resampling

If we do need to create large numbers of external samples via cosmological simulations, Schneider et al. (2011) have proposed a method to rapidly generate, pseudo-independent random realisations from a single N-body simulation. This resampling the large-scale, quasi-Gaussian Fourier modes using a semi-analytic formalism. The approach reproduces power spectrum covariances well, including the coupling between linear and non-linear scales, although a small bias in the covariance elements is introduced. Schneider et al. (2011) find that the number of full N-body simulations required to

achieve the same error tolerance on the covariance matrix is reduced by a factor of 8, at the price of having to run these simulations into the future and with more frequent snapshot outputs.

7.4.2 Internal resampling: Jackknife

If the statistical properties of the data are poorly known, one can create internal samples by resampling the observed data itself. In this case the data covariance is only estimated at one point in parameter-space, from a single realisation. A long-established method is the Jackknife method (Tukey, 1958), which, in the astronomical context of a correlated spatial random process, requires the survey to be split up into N_{sub} equally sized sub-regions. Jackknife samples are constructed by deleting one sub-region in turn (the delete-one Jackknife) and using the galaxy catalogues of the remaining survey area to re-compute the signal mean. The Jackknife covariance of this mean is then given by (Efron, 1980)

$$\widehat{M}_{\text{Jack}} = \frac{N_{\text{sub}} - 1}{N_{\text{sub}}} \sum_{\alpha=1}^{N_{\text{sub}}} \Delta D_{\alpha} \Delta D_{\alpha}, \quad (92)$$

with $\Delta D_{\alpha} = D_{\alpha} - \bar{D}_{\alpha}$, where the subscript α indicates both the realisation and that the sub-region α has been deleted. In the limit of uncorrelated data equation (92) is equivalent to the standard estimator, equation (10) applied to the survey sub-regions (Efron, 1980; Shao & Wu, 1989). In this case there is no advantage in using Jackknifing. If there are correlations between sub-regions, these will be missed by standard estimation while the Jackknife, taken over the whole survey bar one sub-region, will measure these.

We can estimate the number of sub-regions required for the Jackknife, by using the results for the Wishart distribution, assuming that correlations between the sub-regions are negligible, where $N_{\text{sub}} = N_S$. For example in the case of a weak lensing survey covering $15,000 \text{ deg}^2$, the requirement of having of order $N_S = 10^4$ realisation implies that the sub-regions would be little more than 1 degree on a side[†]. To avoid bias due to the impact of sub-region boundaries, the largest scales that could be probed would have to be much smaller, and hence jackknife estimates are likely restricted to, but potentially useful on, smaller scales. Even in this limit, the estimator based on equation (92) will still be biased because the sub-regions are not independent, due to residual correlations, so that the actual number of degrees of freedom is smaller but unknown. In this case, we do not know how to correct for the Wishart bias.

7.4.3 Internal resampling: Bootstrap

Another class of resampling techniques, containing a large number of variants, is the Bootstrap (Efron, 1979). This assumes that the empirical distribution function of a sample of independently and identically distributed data of size n

[†] Note that we do not consider to split up the survey into sub-volumes as done in Norberg et al., (2009) because of the very strong correlations expected along the line-of-sight for a weak lensing survey, due to photometric redshift errors and particularly the broad lensing kernel.

provides an unbiased estimate of the true underlying population from which the data has been drawn. Point estimates, e.g. of the covariance, can be obtained via standard estimators such as equation (10), by generating new samples from this empirical distribution, where the sum now runs over the number of bootstrap samples, N_B . If the bootstrapped sample has size N_r , the maximum number of distinct samples that can be drawn is $(n+N_r-1)!/[(n-1)!N_r!]$, which quickly grows large for moderate n and N_r , so that usually Monte Carlo methods are employed to create bootstrap samples. Sub-regions of the survey are resampled, analogous to the jackknife. In astronomy these resampled sub-regions are generally chosen to be non-overlapping (fixed-block bootstrap), although this may be suboptimal (Nordman, et al., 2007).

If N_B is large, the uncertainty on the covariance estimate becomes negligible, so that its inverse is an unbiased estimator of the precision matrix. However, there are multiple other sources of bias in the bootstrap technique which are hard to quantify. If the number of sub-regions is small, the empirical distribution is a coarse representation of the underlying population, and any local features may be missed, and the convergence of the error tolerance is slower (Nordman, et al., 2007). If the number of sub-regions is large and their area small, the correlation structure is not well preserved. It is possible that some sub-regions are not drawn at all in a given bootstrap sample, so that the total area coverage of this sample is less than that of the original survey, biasing the covariance estimate upwards. Drawing more sub-regions for bootstrap samples than the N_{sub} regions of the original data, i.e. $N_r > N_{\text{sub}}$, can remedy this bias (Norberg, et al., 2009), but the choice of N_r constitutes another parameter that requires calibration.

While the bootstrap method effectively by-passes the Wishart bound by creating a very large number of realisations, its multiple sources of bias, which all depend on the data set at hand, make an application to precision measurements questionable. Note that this is an active field of research, so that this conclusion could change on moderate time scales (see e.g. Loh, 2008, for a recent application of a spatial bootstrap variant to large-scale structure data).

7.5 Summary of Alternatives

Which of the various routes to inverse covariance estimation is optimal depends strongly on the problem at hand, being influenced by aspects as diverse as the survey characteristics, the complexity of obtaining the signal, the computational cost of acquiring simulated survey realisations, the availability and accuracy of models, and the questions one aims to answer with the data. Clearly, if one has an accurate model of the data covariance to hand, this is what should be used. A poorer model can still be used, either to aid a fitting-function approach, or as the model target for shrinkage. Even if a model is not available, empirical shrinkage appears to work well, and all of these methods seem unbiased and yield a statistical error on the precision matrix less than the 5% even when $N_S < N_D$. Stein and Haff shrinkage yield a bias which must be calibrated. Data compression may optimally reduce the number of data-points to the number of model parameters, typically a few hundred values, at the expense of loss of some information, and model-dependency. However, if not optimal the gain may only be a factor of 10's.

Mode-resampling of simulations may lead to production of large numbers of realisation, with a factor of ≈ 8 possible. Empirical resampling of the data will lead to bias estimates which will need calibration, but the Jackknife has promise in the small-scale, highly-subsampled regime. Table 1 provides a summary of our findings. Finally, we caution again that our tests have been on idealised data, which closely match our model, while data compression and resampling methods require detailed testing.

8 SUMMARY AND CONCLUSIONS

Over the next decade and beyond the size of cosmological data sets, and the potential accuracy and ability to probe new physics, will continue to rise dramatically. To ensure that the expected accuracies are reached, we need to consider the dominant sources of bias and uncertainty in our measurements. In this paper we have developed and explored a new framework to study the effect of random errors in the estimation of the data covariance matrix and its inverse, the precision matrix, on cosmological parameter estimation. For multivariate Gaussian data, the likelihood function depends sensitively on the precision matrix.

In many areas of cosmology and astrophysics, the data covariance matrix cannot be predicted analytically and we must rely on estimating it from independent, random realisations of the observations. The simplest estimator of the data covariance matrix is unbiased, and the uncertainty drops with the inverse square-root of the number of independent samples. More generally, the sample data covariance matrix follows a Wishart distribution. In contrast, the simplest estimator for the precision matrix, taking the inverse of the sample data covariance, is biased. We can find an unbiased estimate if the difference between number of realisations, N_S , and number of data points N_D , is $N_S - N_D > 2$. However, the precision matrix follows an Inverse-Wishart distribution, and its variance can diverge if $N_S - N_D \leq 4$.

We have tested and illustrated this behaviour by simulating 10^4 weak lensing surveys and dividing into 100 groups of $N_S = 100$ samples, and shown the mean and variance of the sample data covariance and precision matrix follow the predicted properties of the Wishart and Inverse-Wishart distribution. Future cosmological surveys will have of order $N_D \approx 10^4 - 10^6$ data-points, even with radical compression into power spectra and correlation functions, and so the Wishart bound implies large numbers of realisations will be required. In addition, the only known unbiased estimator of the precision matrix is also the simplest, so other methods will be biased and we may not be able to quantify this bias without comparing with simulations.

The properties of the precision matrix where propagated into the uncertainty in the errors on cosmological parameters, using a Fisher matrix formalism, and we have shown that:

- The fractional errors on the variance of a parameter and the diagonals of the precision matrix are equal.
- The fractional error on the parameter variance depends on inverse square-root of $N_S - N_D - 4$, which can diverge when the number of realisations, N_S is equal to $N_D + 4$.
- The number of realisations needed to reach a given accuracy on parameter errors must be greater than the sum

METHOD		Gain	Cost	Comment
Unbiased		Unbiased, known variance	$N_S > 200 + N_D + 4$	Costly for large N_D
Modelling		Unbiased/No variance	Complex modelling	Nonlinear, baryonic physics
Shrinkage	Stein	factor ≈ 5 reduction in N_S	Unknown bias	Range $N_S > N_D$
	Haff	factor ≈ 5 reduction in N_S	Unknown bias	Range $N_S > N_D$
	Model Target	factor ≈ 10 reduction in N_S	Low/no bias	Applicable for $N_S < N_D$
	Mean Target	factor ≈ 5 reduction in N_S	Low/no bias	Applicable for $N_S < N_D$
Compression		$N_D \approx N_{\text{Para}} \sim 100$	Information loss	Model dependent
Resampling	Simulation	Factor ≈ 8 gain	More expensive simulations	Accuracy to be tested
	Jackknife	No simulations required	Unknown bias and variance	Needs calibration
	Bootstrap	No simulations required /No bias	Multiple unknown biases	Needs calibration

Table 1. Summary Table of results of different methods for estimating the precision matrix for parameter estimation. The methods are described in detail in Section 7.2, and the results of the tests in Section 7.2.4.

of the number data points, and the inverse of the fractional variance of the parameter variance (equation 58).

- The error on the sample data covariance is equal to the precision matrix for small data-sets, while it scales as the inverse-square root of the number of data points for large data sets.

If we want to have a 5% accuracy on a parameter error, and the number of data points is $\ll 100$, we need ≈ 200 realisations and a 10% error on the data covariance matrix. If the data set is greater than 100 points, we will need $N_S \approx N_D$ realisations and the fractional error on the data covariance matrix is $\approx 10(N_D/200)^{-1/2}$ percent. We also have shown how the uncertainties propagates into the Figure-of-Merit, and found similar conclusions. To attain high-accuracy from large-scale data sets seems to require equally large-numbers of realisations of the survey. We have explored some of the possible alternatives to alleviate this conclusion:

- **Theoretical Modelling:** If we can accurately model the data covariance we avoid the Wishart bound.
- **Shrinkage Estimators:** We reduce the required sample by combining empirical or theoretical estimates of the precision or data covariance matrix with sample estimates.
- **Data Compression:** We can reduce the number of data points, N_D , in principle to the number of parameters.
- **Simulations mode resampling:** We can rapidly generate external realisations, but should check for accuracy.
- **Data Resampling:** The Jackknife resampling method may be useful on small-scales but will be biased. The Bootstrap method can generate large number of samples and so not be biased, but needs further study to avoid other biases.

The combination of theoretical modelling and target shrinkage looks particularly promising and robust, but clearly needs to be developed and tested in more detail for application.

In summary, many of the details of precision cosmology are still to be worked out. We have identified the estimation of the precision matrix as a key issue, for Gaussian-distributed data, requiring the generation of large numbers of realisations for large-datasets. We have investigated a

number of possible ways forward, although the actual resolution of this issue may require the use of multiple techniques.

ACKNOWLEDGEMENTS

We thank Martin White for useful discussion about the effect of the precision matrix on the parameter error, Peder Norberg for useful discussion about the Jackknife and Bootstrap methods, and Alina Kiessling for encouraging our interest in this problem. We also thank an anonymous referee of useful comments. BJ thanks the STFC for funding on a Consolidated Grant, while TDK acknowledges the support of a Royal Society University Research Fellowship.

REFERENCES

- Albrecht A., et al. (The Dark Energy Task Force), 2006 (arXiv:astro-ph/0609591)
- Amendola L., et al. (Euclid Theory Working Group), 2012 (arXiv:1206.1225)
- Anderson, T. W. 2003, An introduction to multivariate statistical analysis, 3rd edn. (Wiley-Interscience)
- Asgari M., Schneider P., Simon P., 2012, submitted (arXiv:1201.2669)
- Ballinger W.E., Heavens A.F., Taylor A.N., MNRAS, 276, 59
- Bird S., Viel M., Haehnelt M.G., , 2012, MNRAS, 420, 2551
- Bond J.R., Jaffe A.H., Knox L., 1998, Phys.Rev. D, 57, 2117
- Brown M., Taylor A.N., Bacon D.J., Gray M.E., Dye S., Meisenheimer K., Wolf C., 2003, MNRAS, 341, 100
- Brown M.L., Castro P.G., Taylor A.N., 2005, MNRAS, 360, 1262
- Bucher M., Moodley K., Turok N., 2001, Phys. Rev. Lett., 87, 191301
- Casaponsa B., Heavens A.F., Kitching T.D., Miller L., Barreiro R.B., Martinez-Gonzalez E., 2012, MNRAS, submitted (arXiv:1209.1646)
- Chevallier M., Polarski D., 2001, IJMPD, 10, 213
- Clifton T., Ferreira P.G., Padillo A., Skordis C., 2012, Phys. Rep., 513, 1
- Cooray A., Hu W., 2001, ApJ, 554, 56
- Copeland E., Sami M., Tsujikawa S., 2006, Int. J. Mod. Phys. D, 15, 1753
- Efron B., 1979, Ann. Statist., 7, 1

- Efron B., 1980, The Jackknife, the Bootstrap and Other Resampling Plans, CBMS-NSF Regional Conference Series in Applied Mathematics
- Eriksen H.K., Wehus I.K., 2009, *ApJ Supp.*, 180, 30
- Gupta, S., Heavens A.F., 2002, *MNRAS*, 334, 167
- Haff, L.R., 1974, *Ann. Statist.*, 7, 1264
- Hamilton A.J.S., 1998, in *The Evolving Universe: Selected Topics on Large-Scale Structure and on the Properties of Galaxies*, Dordrecht: Kluwer Academic Publishers
- Hamilton A.J.S., Rimes C.D., Scoccimarro R., 2005, *MNRAS*, 371, 1188
- Hamimeche S., Lewis A., 2009, *Phys. Rev. D*, 79, 083012
- Hartlap J., Simon P., Schneider P., 2007, *A&A*, 464, 399
- Heavens A.F., Taylor A.N., 1995, *MNRAS*, 275, 483
- Heavens A.F., Jimenez R., Lahav O., 2000, *MNRAS*, 317, 965
- Heavens A.F., 2003, *MNRAS*, 343, 1327
- Hilbert S., Hartlap J., Schneider P., 2011, *A&A*, 536, 85
- Hu W., 1999, *ApJ*, 522, 21
- Huterer D., Takada M., 2005, *Astropart. Phys.*, 23, 369
- Huterer D., Takada M., Bernstein G., Jain B., 2006, *MNRAS*, 366, 101
- Kaiser N., 1987, *MNRAS*, 227, 1
- Kaiser N., 1992, *ApJ*, 388, 272
- Kaufman G.M., 1967, Some Bayesian Moment Formulae, Report No. 6710, Center for Operations Research and Econometrics, Catholic University of Louvain, Heverlee, Belgium
- Kayo I., Takada M., Jain B., 2012, submitted (arXiv:1207.6322)
- Kiessling A., Taylor A.N., Heavens A.F., 2011, *MNRAS*, 416, 1045
- Kitching T.D., Heavens A.F., Taylor A.N., Brown M.L., Meisenheimer K., Wolf C., Gray M.E., Bacon D.J., 2007, *MNRAS*, 374, 771
- Knox L., 1995, *Phys. Rev. D*, 58, 123506
- Larson D.L., Eriksen H.K., Wandelt B.D., Gorski K.M., Huey G., Jewell J.B., O'Dwyer I.J., 2007, *ApJ*, 656, 653
- Ledoit O., Wolf M., 2003, *Journal of Empirical Finance*, 10, 603
- Liddle A., Mukherjee P., Parkinson D., 2006, *A&G*, 47, 4.30
- Linder E.V., 2003, *Phys. Rev. Lett.*, 90, 091301
- Loh J.M., 2008, *ApJ*, 681, 726
- Matsumoto S., 2011, *Journal of Theoretical Probability*, Online First, 1 (arXiv:1004.4717v3 [math.ST])
- Meiksin A., White M., 1999, *MNRAS*, 308, 1179
- Norberg P., Baugh C.M., Gaztañaga E., Croton D.J., 2009, *MNRAS*, 396, 19
- Nordman D.J., Lahiri S.N., Fridley B.L., 2007, *Indian Journal of Statistics*, 69, 468
- Percival W.J., Brown M.L., 2006, *MNRAS*, 372, 1104
- Pope A.C., Szapudi I., 2008, *MNRAS*, 389, 766
- Press, S.J., 1982, *Applied Multivariate Analysis*, Robert E. Krieger Publishing Company, Florida
- Rimes C.D., Hamilton A.J.S., 2005, *MNRAS*, 360, L82
- Rimes C.D., Hamilton A.J.S., 2006, *MNRAS*, 371, 1205
- Sato M., Hamana T., Takahashi R., Takada M., Yoshida N., Matsubara T., Sugiyama N., 2009, *ApJ*, 701, 945
- Schäfer J., Strimmer K., 2005, *Statist. App. Mol. Genet. Biol.*, 4, A32
- Schmidt F., Leauthaud A., Massey R., Rhodes J., George M., Koekemoer A.M., Finoguenov A., Tanaka M., 2012, *ApJ*, 744, L22
- Schneider M., Cole S., Frenk C., Szapudi I., 2011, *ApJ*, 737, 11
- Schneider P., Eifler T., Krause E., 2010, *A&A*, 520, 116
- Scoccimarro R., Zaldarriaga M., Hui L., 1999, *ApJ*, 527, 1
- Shao J., Wu C.F.J., 1989, *Ann. Statist.*, 17, 1176
- Sivia, D.S., 1996, *Data Analysis: A Bayesian Tutorial*, Clarendon Press, Oxford
- Spergel D. N., Verde L., Peiris H. V., Komatsu E., Nolte M. R., Bennett C. L., Halpern M., Hinshaw G., Jarosik N., Kogut A., Limon M., Meyer S. S., Page L., Tucker G. S., Weiland J. L., Wollack E., Wright E. L., *ApJSupp*, 148, 175
- Stein C., Efron B., Morris C., 1972, Tech. Report No. 37, Depart. Statist., Stanford Univ.
- Tadros, H., et al., 1999, *MNRAS*, 305, 527
- Takada M., Bridle S., 2007, *New Journal of Physics*, 9, 446
- Takada M., Jain B., 2009, *MNRAS*, 395, 2065
- Takahashi R.m Yoshida N., Takada M., Matsubara T., Sugiyama N. Kayo I., Nishizawa A.J., Nishimichi T., Saito S., Taruya A., 2009, *ApJ*, 700, 479
- Taylor A.N., Ballinger W.E., Heavens A.F., Tadros H., 2001, *MNRAS*, 327, 689
- Taylor A.N., Kitching T.D., Bacon D.J., Heavens A.F., 2007, *MNRAS*, 374, 1377
- Taylor A.N., Kitching T.D., 2010, *MNRAS*, 408, 865
- Tegmark M., Taylor A. N., Heavens A. F., 1997, *ApJ*, 480, 22
- Trotta R., 2007, *MNRAS*, 378, 72
- Tukey J., 1958, *Ann. Math. Statist.*, 29, 614
- Verde L., Peiris H. V., Spergel D. N., Nolte M. R., Bennett C. L., Halpern M., Hinshaw G., Jarosik N., Kogut A., Limon M., Meyer S. S., Page L., Tucker G. S., Wollack E., Wright E. L., *ApJSupp*, 148, 195
- Wishart J., 1928, *Biometrika*, 20A, 32

APPENDIX A: COSMOLOGICAL LARGE-SCALE STRUCTURE DATA-SETS

In general, cosmological data can be in any number of forms, and our analysis is applicable to a wide variety of data. If we are working with pixelised maps then the data are pixel-values and the data covariance matrix is the pixel-covariance matrix. If we have compressed information into two-point power spectra or correlation functions then this forms the data vector, and the data covariance matrix is the field's four-point function. In this paper we shall assume the data is compressed into the two-point power spectra, although our formulae can be used for pixelised data, correlation functions, or higher-order correlations. Here we outline three basic areas of large-scale structure study, and how their data-vectors are generated.

A1: Galaxy Redshift Surveys

In galaxy redshift surveys, we can predict the distribution of the matter overdensity field,

$$\delta(\mathbf{r}) = \frac{\rho(\mathbf{r}) - \langle \rho \rangle}{\langle \rho \rangle}, \quad (93)$$

which we can compare with data about the galaxy overdensity,

$$\delta_g(\mathbf{r}) = \frac{n(\mathbf{r}) - \bar{n}(r)}{\bar{n}(r)}, \quad (94)$$

where $n(\mathbf{r})$ is the galaxy number-distribution and $\bar{n}(r)$ is the survey radial selection function. The Fourier transform of the matter overdensity, $\delta(\mathbf{k})$, is

$$\delta(\mathbf{k}) = \int d^3r \delta(\mathbf{r}) e^{-i\mathbf{k} \cdot \mathbf{r}}. \quad (95)$$

The galaxy distribution can also be expanded in spherical harmonics, $\delta_{\ell m}(z)$, and radial Bessel functions, $\delta_{\ell mn}$ (e.g., Heavens & Taylor, 1995). As we measure galaxy radial positions with redshift, which combines the Hubble expansion with peculiar velocities, the galaxy distribution is changed

by redshift-space distortions (Kaiser, 1987; Hamilton, 1998) so that in Fourier space,

$$\delta_g^s(\mathbf{k}) = (b(k) + f(\Omega_m)\mu_k^2)\delta(\mathbf{k}), \quad (96)$$

where $b(k)$ is a scale-dependent galaxy bias factor, $f(\Omega_m) = d\ln\delta/d\ln a$ is the growth index. The correlation of the modes of the overdensity field is

$$\langle \delta_g^s(\mathbf{k}) \delta_g^{s*}(\mathbf{k}') \rangle = (2\pi)^3 P_{gg}^s(\mathbf{k}) \delta_D(\mathbf{k} - \mathbf{k}'), \quad (97)$$

where $P_{gg}^s(\mathbf{k})$ is the anisotropic redshift-space galaxy power spectrum. We can estimate the anisotropic redshift-space galaxy density power spectrum from

$$\hat{P}_{gg}^s(k, \mu_k) = \frac{1}{N_{\text{modes}}} \sum_{k_z} |\delta_g^s(\mathbf{k})|^2, \quad (98)$$

where the summation is over the N_{modes} in the k_z -direction. The data is the discretely sampled redshift-space power spectrum;

$$D_i = \hat{P}_{gg}^s(\mathbf{k}_i), \quad (99)$$

where a hat $\hat{}$ indicates the observed estimate of the power.

A2: CMB

In Cosmic Microwave Background experiments we can define the temperature fluctuations as $\Theta = \Delta T/T$, and the temperature power-spectrum, C_ℓ^{TT} , is defined by

$$\langle \Theta_{\ell m} \Theta_{\ell' m'}^* \rangle = C_\ell^{TT} \delta_{\ell\ell'} \delta_{mm'}. \quad (100)$$

If polarisation data is added to this, in the form of E - and B -modes, we can construct 6 power spectra,

$$\mathbf{D} = (\hat{C}_\ell^{TT}, \hat{C}_\ell^{EE}, \hat{C}_\ell^{BB}, \hat{C}_\ell^{TE}, \hat{C}_\ell^{TB}, \hat{C}_\ell^{EB}), \quad (101)$$

which can form our data-vector. We can estimate these cross-spectra from

$$\hat{C}_\ell^{XY} = \frac{1}{(2\ell+1)} \sum_m X_{\ell,m} Y_{\ell,m}^*, \quad (102)$$

where $(X, Y) = (\Theta, E, B)$, and we have summed over all azimuthal m -modes for each ℓ . Again, in practise these power-spectra would be convolved by the survey mask (e.g., Brown, Castro & Taylor, 2005).

A3: Weak Lensing

In weak lensing surveys the data can be the estimated shear values, $\gamma_i(\boldsymbol{\theta}, z)$, where $i = (1, 2)$ are the two orthogonal modes of the shear. The shear-shear covariance matrix, for an unmasked survey, is

$$\langle \gamma_i(\boldsymbol{\ell}, z) \gamma_j^*(\boldsymbol{\ell}', z') \rangle = (2\pi)^2 C_{ij}^{\gamma\gamma}(\boldsymbol{\ell}, z, z') \delta_D(\boldsymbol{\ell} - \boldsymbol{\ell}') \quad (103)$$

where $C_{ij}^{\gamma\gamma}(\boldsymbol{\ell}, z, z')$ is the shear power-spectrum. The shear can be decomposed into a potential (κ , convergence) and curl (β) part,

$$\kappa(\boldsymbol{\ell}, z) + i\beta(\boldsymbol{\ell}, z) = e^{2i\varphi_\ell} (\gamma_1 + i\gamma_2)(\boldsymbol{\ell}, z), \quad (104)$$

where φ_ℓ is the angle between the wavevector, $\boldsymbol{\ell}$, and an axis of the coordinate system the shear is measured in. This decomposition generates three power spectra, $C^{\kappa\kappa}(\boldsymbol{\ell}, z, z')$, $C^{\beta\beta}(\boldsymbol{\ell}, z, z')$, $C^{\kappa\beta}(\boldsymbol{\ell}, z, z')$. In principle, we can also add a magnification field, μ , estimated from the size of

galaxy images (e.g., Schmidt et al., 2012; Casaponsa et al., 2012), yielding the magnification power, $C^{\mu\mu}(\boldsymbol{\ell}, z, z')$. The data, compressed into these power spectra, is then

$$\mathbf{D} = (\hat{C}^{\mu\mu}(\boldsymbol{\ell}, z, z'), \hat{C}^{\kappa\kappa}(\boldsymbol{\ell}, z, z'), \hat{C}^{\beta\beta}(\boldsymbol{\ell}, z, z'), \hat{C}^{\mu\kappa}(\boldsymbol{\ell}, z, z'), \hat{C}^{\mu\beta}(\boldsymbol{\ell}, z, z'), \hat{C}^{\kappa\beta}(\boldsymbol{\ell}, z, z')). \quad (105)$$

The estimated cross-power spectrum, at two different redshifts, can be estimated by

$$\hat{C}_\ell^{XY}(z, z') = \frac{1}{N_{\text{modes}}} \sum_{|\boldsymbol{\ell}|=\ell} X(\boldsymbol{\ell}, z) Y^*(\boldsymbol{\ell}, z'), \quad (106)$$

where $(X, Y) = (\mu, \kappa, \beta)$, and we have summed over N_{modes} is a shell in ℓ -space. In general, the effects of a survey mask, due to survey geometry and bright stars, will convolve these spectra. We can also estimate these statistics in real-space, through correlation functions, or other weighted two-point functions such as M_{ap} or COSEBIs (Schneider, et al., 2010; Asgari, et al., 2012).

A4: Combining data-sets

Combining data-sets can be achieved by combining the data-vectors of each data-set, and including all the cross-terms between the surveys. To do this we define a vector for field-values;

$$\mathbf{X} = (\delta_{g,\ell m}(z), \Theta_{\ell m}, E_{\ell m}, B_{\ell m}, \mu_{\ell m}(z), \kappa_{\ell m}(z), \beta_{\ell m}(z)), \quad (107)$$

where we have chosen to expand the lensing shear, magnification and galaxy redshift fields into spherical harmonics for consistency. We then form all of the auto- and cross-spectra of these fields,

$$C_\ell^{X_i X_j}(z, z') = \frac{1}{2\ell+1} \sum_m X_i(\boldsymbol{\ell}, z) X_j^*(\boldsymbol{\ell}, z'), \quad (108)$$

where we have summed over all azimuthal modes on the sky. The data-vector, \mathbf{D} , is then the vector of all auto- and cross-spectra.

APPENDIX B: PROPERTIES OF THE WISHART AND INVERSE-WISHART DISTRIBUTIONS

Many of the properties of the Wishart and Inverse-Wishart distributions reside in technical mathematical statistics papers and specialist textbooks. Few proofs written for physicists are available, so in this Appendix we present our own derivations of some useful results used in this paper. Most of the results, if not the details, can be found in Press (1982).

B1. Wishart distribution

Let \mathbf{V} be a $p \times p$ symmetric matrix so that

$$\mathbf{V} = \sum_{\alpha=1}^n \mathbf{x}_\alpha \mathbf{x}_\alpha^t, \quad (109)$$

where \mathbf{x} is a vector drawn from a multivariate Gaussian distribution, and each of n realisations is sampled independently. The distribution of \mathbf{V} is given by (Wishart, 1928)

$$p(\mathbf{V}|\mathbf{\Sigma}) = c|\mathbf{V}|^{(n-p-1)/2}|\mathbf{\Sigma}|^{-n/2}e^{-\frac{1}{2}\text{Tr}\mathbf{V}\mathbf{\Sigma}^{-1}}, \quad (110)$$

where

$$c = [2^{np/2}\Gamma_p[n/2]]^{-1}, \quad (111)$$

and $\mathbf{\Sigma}$ sets the scale of the distribution. The Multivariate Gamma Function, $\Gamma_p(a)$, is defined as

$$\Gamma_p(a) \equiv \int_{\mathbf{X}>0} d\mathbf{X} |\mathbf{X}|^{a-(p+1)/2} e^{-\text{Tr}\mathbf{X}} \quad (112)$$

$$= \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma\left(a + \frac{1-j}{2}\right), \quad (113)$$

where \mathbf{X} is a $p \times p$ matrix and the matrix integration is over all positive-definite elements,

$$d\mathbf{X} \equiv \prod_{i=1}^p \prod_{j=1}^p dX_{ij}, \quad (114)$$

or for a symmetric matrix over all non-repeated elements,

$$d\mathbf{X} \equiv \prod_{i=1}^p \prod_{j=1}^i dX_{ij}. \quad (115)$$

B2. Derivation of the Wishart distribution

We can derive the Wishart distribution through its characteristic function, the Fourier transform of the probability distribution function,

$$\phi(\mathbf{J}) = \int d\mathbf{V} p(\mathbf{V}|\mathbf{\Sigma}) e^{i\text{Tr}\mathbf{J}\mathbf{V}} = \langle e^{i\text{Tr}\mathbf{J}\mathbf{V}} \rangle. \quad (116)$$

Expanding the Fourier exponential in a Taylor series and taking expectations we find,

$$\phi(\mathbf{J}) = 1 + i\langle\text{Tr}(\mathbf{J}\mathbf{V})\rangle - \frac{1}{2}\langle(\text{Tr}\mathbf{J}\mathbf{V})(\text{Tr}\mathbf{J}\mathbf{V})\rangle + \dots \quad (117)$$

Using the Gaussian properties of \mathbf{x} , the first and second moments of \mathbf{V} are

$$\langle\mathbf{V}\rangle = \sum_{\alpha=1}^n \langle\mathbf{x}_{\alpha}\mathbf{x}_{\alpha}^t\rangle = n\mathbf{\Sigma}, \quad (118)$$

$$\begin{aligned} \langle V_{ij}V_{mn} \rangle &= \left\langle \sum_{\alpha=1}^n x_{i,\alpha}x_{j,\alpha} \sum_{\beta=1}^n x_{m,\beta}x_{n,\beta} \right\rangle \\ &= n(\Sigma_{im}\Sigma_{jn} + \Sigma_{in}\Sigma_{jm}) + n^2\Sigma_{ij}\Sigma_{mn}. \end{aligned} \quad (119)$$

Taking the expectation values, we can re-write the series for the characteristic function in terms of the scale matrix, $\mathbf{\Sigma}$, as

$$\begin{aligned} \phi(\mathbf{J}) &= 1 + i n \text{Tr}(\mathbf{J}\mathbf{\Sigma}) \\ &\quad - \frac{1}{2} (2n \text{Tr}(\mathbf{J}\mathbf{\Sigma}\mathbf{J}\mathbf{\Sigma}) + n^2 (\text{Tr}\mathbf{J}\mathbf{\Sigma})(\text{Tr}\mathbf{J}\mathbf{\Sigma})) \\ &\quad + \dots \end{aligned} \quad (120)$$

Collecting terms in equation (120) in powers of n we find

$$\begin{aligned} \phi(\mathbf{J}) &= 1 + n[i\text{Tr}(\mathbf{J}\mathbf{\Sigma}) - \text{Tr}[(\mathbf{J}\mathbf{\Sigma})^2] + \dots] \\ &\quad + \frac{1}{2}n^2[i\text{Tr}(\mathbf{J}\mathbf{\Sigma}) - \text{Tr}[(\mathbf{J}\mathbf{\Sigma})^2] + \dots]^2 + \dots \end{aligned} \quad (121)$$

Each of the series with factor n , n^2 , etc, can be summed to a logarithm, using the series relation $\ln(1+x) =$

$\sum_{n=0}^{\infty} (-1)^{n+1} x^n / n$, which yields

$$\begin{aligned} \phi(\mathbf{J}) &= 1 - \frac{n}{2} \text{Tr} \ln(\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}) + \frac{n^2}{8} [\text{Tr} \ln(\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma})]^2 \\ &\quad + \dots \end{aligned} \quad (122)$$

This last series can now be summed, using $e^x = \sum_{n=0}^{\infty} x^n / n!$, to find

$$\phi(\mathbf{J}) = \exp\left(-\frac{n}{2} \text{Tr} \ln(\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma})\right). \quad (123)$$

Using the matrix identity, $\ln \det \mathbf{A} = \text{Tr} \ln \mathbf{A}$, we find the characteristic function for \mathbf{V} is

$$\phi(\mathbf{J}) = |\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}|^{-n/2}. \quad (124)$$

We can show the Wishart distribution has the same characteristic function by direct integration;

$$\begin{aligned} \phi(\mathbf{J}) &= \int d\mathbf{V} p(\mathbf{V}|\mathbf{\Sigma}) e^{i\text{Tr}\mathbf{J}\mathbf{V}} \\ &= c|\mathbf{\Sigma}|^{-n/2} \int d\mathbf{V} |\mathbf{V}|^{(n-p-1)/2} e^{-\frac{1}{2}\text{Tr}\mathbf{V}\mathbf{\Sigma}^{-1}[\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}]} \end{aligned} \quad (125)$$

It is convenient to define the new matrix variable

$$\mathbf{\Theta} = \mathbf{\Sigma}^{-1}[\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}], \quad (126)$$

and the matrix

$$\mathbf{X} = \frac{1}{2}\mathbf{V}\mathbf{\Theta}. \quad (127)$$

The transformation of the matrix volume element is given by

$$d\mathbf{V} = 2^{p(p+1)/2} |\mathbf{\Theta}|^{-(p+1)/2} d\mathbf{X}. \quad (128)$$

We can now rewrite the characteristic function as

$$\begin{aligned} \phi(\mathbf{J}) &= c|\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}|^{-n/2} 2^{pn/2} \int d\mathbf{X} |\mathbf{X}|^{(n-p-1)/2} e^{-\text{Tr}\mathbf{X}} \\ &= |\mathbf{I} - 2i\mathbf{J}\mathbf{\Sigma}|^{-n/2}, \end{aligned} \quad (129)$$

where we have made use of the Multivariate Gamma Function (equation 113) to cancel terms in c . Identifying $\phi(\mathbf{J})$ as the characteristic function of \mathbf{V} from equation (124), we confirm that the Wishart distribution, $p(\mathbf{V}|\mathbf{\Sigma})$, is the probability distribution for the matrix \mathbf{V} . The moments of the Wishart distribution can be found directly from the expansion of the characteristic function in equation (120).

B3. Inverse-Wishart distribution

The Inverse-Wishart distribution may be found from the Wishart distribution by a change of variables. We first note that the Jacobian for the transformation $\mathbf{U} = \mathbf{V}^{-1}$, where \mathbf{U} and \mathbf{V} are $p \times p$ symmetric matrices is (see Appendix C3)

$$d\mathbf{V} = |\mathbf{U}|^{-(p+1)} d\mathbf{U}. \quad (130)$$

If we further define $\mathbf{G} = \mathbf{\Sigma}^{-1}$ then

$$\begin{aligned} p(\mathbf{U}|\mathbf{G})d\mathbf{U} &= p(\mathbf{V}|\mathbf{\Sigma})d\mathbf{V} \\ &= c|\mathbf{V}|^{(n-p-1)/2} |\mathbf{\Sigma}|^{-n/2} e^{-\frac{1}{2}\text{Tr}\mathbf{V}\mathbf{\Sigma}^{-1}} d\mathbf{V}, \\ &= \left[c|\mathbf{U}|^{-(n-p-1)/2} |\mathbf{G}|^{n/2} e^{-\frac{1}{2}\text{Tr}\mathbf{U}^{-1}\mathbf{G}} \right] \\ &\quad \times |\mathbf{U}|^{-(p+1)} d\mathbf{U}, \end{aligned}$$

$$= c |\mathbf{U}|^{-(n+p+1)/2} |\mathbf{G}|^{n/2} e^{-\frac{1}{2} \text{Tr} \mathbf{U}^{-1} \mathbf{G}} d\mathbf{U} \quad (131)$$

Hence, the Inverse-Wishart distribution is

$$p(\mathbf{U}|\mathbf{G}) = c |\mathbf{U}|^{-(n+p+1)/2} |\mathbf{G}|^{n/2} e^{-\frac{1}{2} \text{Tr} \mathbf{U}^{-1} \mathbf{G}}. \quad (132)$$

We should note that we have assumed the number-of-degrees of freedom, n , is the same for the Wishart and Inverse-Wishart, to simplify the derivation. However, the Inverse-Wishart can be parameterized differently, and different authors choose different relations between the Wishart and Inverse-Wishart degrees-of-freedom. If we let m be the number of degrees-of-freedom for the Inverse-Wishart, we have set $m = n$. Other choices commonly used are $m = n + p - 1$, $m = n + p + 1$, or $m = n - p - 1$. For example if we assume $m = n - p - 1$ (Press, 1982), we would write the Inverse-Wishart distribution as

$$p(\mathbf{U}|\mathbf{G}) = c_0 |\mathbf{G}|^{(n-p-1)/2} |\mathbf{U}|^{-n/2} e^{-\frac{1}{2} \text{Tr} \mathbf{G} \mathbf{U}^{-1}}, \quad (133)$$

where

$$c_0 = [2^{p(n-p-1)/2} \Gamma_p[(n-p-1)/2]]^{-1}. \quad (134)$$

The moments of the Inverse-Wishart can be found by direct integration over the Inverse-Wishart distribution.

APPENDIX C: CHANGE OF RANDOM MATRIX VARIABLES

In many cases we want to know how to change random matrix variables, for example in order to carry out matrix integration and to derive the Inverse-Wishart distribution. Here we present some useful matrix transformation relations, without proof.

C1: Vector transformations

We first consider vector integration and change of variable. For a p -dimensional vector \mathbf{x} , the infinitesimal volume-element is

$$d\mathbf{x} = d^p x = \prod_{i=1}^p dx_i. \quad (135)$$

For a vector, \mathbf{x} which is related to the vector \mathbf{y} by the linear transformation,

$$\mathbf{y} = \mathbf{A}\mathbf{x}, \quad (136)$$

or with indices $y_i = A_{ij}x_j$, and the volume element transforms as

$$\prod_{i=1}^p dy_i = |\mathbf{A}| \prod_{j=1}^p dx_j \quad (137)$$

or equivalently

$$d\mathbf{y} = |\mathbf{A}| d\mathbf{x}. \quad (138)$$

C2: Transformation of non-symmetric matrices

If we transform a $p \times p$ matrix \mathbf{X} to the matrix \mathbf{Y} ,

$$\mathbf{Y} = \mathbf{A}\mathbf{X}, \quad (139)$$

or with indices, $Y_{ij} = A_{ik}X_{kj}$ the matrix volume-elements for a non-symmetric matrix are

$$d\mathbf{X} = \prod_{i=1}^p \left[\prod_{j=1}^p dX_{ij} \right]. \quad (140)$$

Each of the sub-vectors transforms as a vector, so that

$$\prod_{j=1}^p dY_{1j} = |\mathbf{A}| \prod_{j=1}^p dX_{1j}, \quad (141)$$

as each sub-vector has the same determinant of \mathbf{A} . The matrix volume-elements then transform as

$$d\mathbf{Y} = |\mathbf{A}|^p d\mathbf{X}. \quad (142)$$

If we consider now the linear transformation

$$\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{B}, \quad (143)$$

or $Y_{ij} = A_{ik}X_{kl}B_{lj}$, where \mathbf{X} is $p \times q$ and \mathbf{B} is $q \times q$, we can apply our transformation rules twice to see that

$$d\mathbf{Y} = |\mathbf{A}|^p |\mathbf{B}|^q d\mathbf{X}. \quad (144)$$

If $\mathbf{B} = \mathbf{A}^t$, and $q = p$ we find for a general matrix

$$d\mathbf{Y} = |\mathbf{A}|^{2p} d\mathbf{X}. \quad (145)$$

C3: Transformation of symmetric matrices

If the matrix \mathbf{X} is a $p \times p$ symmetric matrix, $\mathbf{X} = \mathbf{X}^t$, and we consider a transformation of the form

$$\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{A}^t \quad (146)$$

then

$$d\mathbf{Y} = |\mathbf{A}|^{p+1} d\mathbf{X}, \quad (147)$$

when $|\mathbf{A}| \neq 0$. In addition, if \mathbf{X} is symmetric and we multiply by a scalar, a , so that

$$\mathbf{Y} = a\mathbf{X}, \quad (148)$$

then

$$d\mathbf{Y} = a^{p(p+1)/2} d\mathbf{X}. \quad (149)$$

Finally, if we want to transform to the inverse of a symmetric matrix, so that

$$\mathbf{Y} = \mathbf{X}^{-1} = \mathbf{X}^{-1} \mathbf{X} \mathbf{X}^{-1}, \quad (150)$$

then

$$d\mathbf{Y} = |\mathbf{X}|^{-(p+1)} d\mathbf{X}. \quad (151)$$